# pandas: powerful Python data analysis toolkit

Release 1.0.5

**Wes McKinney and the Pandas Development Team** 

## **CONTENTS**

1	Getti	ing starte	ed 3
	1.1	Installat	tion
	1.2	Intro to	pandas
	1.3	Coming	g from 5
	1.4	Commu	nnity tutorials
		1.4.1	Installation
		1.4.2	Package overview
		1.4.3	10 minutes to pandas
		1.4.4	Getting started tutorials
		1.4.5	Essential basic functionality
		1.4.6	Intro to data structures
		1.4.7	Comparison with other tools
		1.4.8	Tutorials
2	User	Guide	225
	2.1	IO tools	s (text, CSV, HDF5,)
		2.1.1	CSV & text files
		2.1.2	JSON
		2.1.3	HTML
		2.1.4	Excel files
		2.1.5	OpenDocument Spreadsheets
		2.1.6	Binary Excel (.xlsb) files
		2.1.7	Clipboard
		2.1.8	Pickling
		2.1.9	msgpack
		2.1.10	HDF5 (PyTables)
		2.1.11	Feather
		2.1.12	Parquet
		2.1.13	ORC
		2.1.14	SQL queries
		2.1.15	Google BigQuery
		2.1.16	Stata format
		2.1.17	SAS formats
		2.1.18	SPSS formats
		2.1.19	Other file formats
		2.1.20	Performance considerations
	2.2		g and selecting data
		2.2.1	Different choices for indexing
		2.2.2	Basics
		2.2.3	Attribute access
		2.2.9	Titalione decess

	2.2.4	Slicing ranges	347
	2.2.5	Selection by label	348
	2.2.6	Selection by position	352
	2.2.7	Selection by callable	355
	2.2.8	IX indexer is deprecated	357
	2.2.9	Indexing with list with missing labels is deprecated	358
	2.2.10	Selecting random samples	360
	2.2.11	Setting with enlargement	362
	2.2.12	Fast scalar value getting and setting	363
	2.2.13	Boolean indexing	364
	2.2.14	Indexing with isin	366
	2.2.15	The where () Method and Masking	368
	2.2.16	The query () Method	372
	2.2.17	Duplicate data	382
	2.2.18	Dictionary-like get () method	385
	2.2.19	The lookup () method	
	2.2.20	Index objects	385
	2.2.21	Set / reset index	389
	2.2.22	Returning a view versus a copy	
2.3	MultiIn	dex / advanced indexing	
	2.3.1	Hierarchical indexing (MultiIndex)	395
	2.3.2	Advanced indexing with hierarchical index	
	2.3.3	Sorting a MultiIndex	111
	2.3.4	Take methods	114
	2.3.5	Index types	
	2.3.6	Miscellaneous indexing FAQ	
2.4	Merge,	join, and concatenate	
	2.4.1	Concatenating objects	130
	2.4.2	Database-style DataFrame or named Series joining/merging	
	2.4.3	Timeseries friendly merging	<del>1</del> 61
2.5		ing and pivot tables	164
	2.5.1	Reshaping by pivoting DataFrame objects	
	2.5.2	Reshaping by stacking and unstacking	
	2.5.3	Reshaping by Melt	174
	2.5.4	Combining with stats and GroupBy	
	2.5.5	Pivot tables	
	2.5.6	Cross tabulations	
	2.5.7	Tiling	
	2.5.8	Computing indicator / dummy variables	
	2.5.9	$\epsilon$	486
	2.5.10	1	486
	2.5.11		189
2.6			491
	2.6.1	V 1	491
	2.6.2	E	193
	2.6.3		496
	2.6.4		500
	2.6.5	e	505
	2.6.6		506
	2.6.7		510
	2.6.8	e	511
2.7	2.6.9	· · · · · · · · · · · · · · · · · · ·	511
2.7		g with missing data	
	2.7.1	Values considered "missing"	)13

	2.7.2	Inserting missing data
	2.7.3	Calculations with missing data
	2.7.4	Sum/prod of empties/nans
	2.7.5	NA values in GroupBy
	2.7.6	Filling missing values: fillna
	2.7.7	Filling with a PandasObject
	2.7.8	Dropping axis labels with missing data: dropna
	2.7.9	Interpolation
	2.7.10	Replacing generic values
	2.7.11	String/regular expression replacement
	2.7.12	Numeric replacement
	2.7.13	Experimental NA scalar to denote missing values
2.8	Categor	ical data
	2.8.1	Object creation
	2.8.2	CategoricalDtype
	2.8.3	Description
	2.8.4	Working with categories
	2.8.5	Sorting and order
	2.8.6	Comparisons
	2.8.7	Operations
	2.8.8	Data munging
	2.8.9	Getting data in/out
	2.8.10	Missing data
	2.8.11	Differences to R's factor
		Gotchas
2.9		e integer data type
2.9	2.9.1	Construction
	2.9.1	
		Operations
2 10	2.9.3	Scalar NA Value
2.10		e Boolean Data Type
	2.10.1	Indexing with NA values
		Kleene Logical Operations
2.11		zation
	2.11.1	Basic plotting: plot
		Other plots
	2.11.3	Plotting with missing data
		Plotting Tools
	2.11.5	Plot Formatting
		Plotting directly with matplotlib
2.12	Comput	tational tools
	2.12.1	Statistical functions
	2.12.2	Window Functions
	2.12.3	Aggregation
	2.12.4	Expanding windows
	2.12.5	Exponentially weighted windows
2.13	Group I	By: split-apply-combine
	2.13.1	Splitting an object into groups
	2.13.2	Iterating through groups
	2.13.3	Selecting a group
	2.13.4	Aggregation
	2.13.5	Transformation
	2.13.6	Filtration
	2.13.7	Dispatching to instance methods
	2.13.8	Flexible apply

	2.13.9	Other useful features
	2.13.10	Examples
2.14	Time se	ries / date functionality
	2.14.1	Overview
	2.14.2	Timestamps vs. Time Spans
	2.14.3	Converting to timestamps
	2.14.4	Generating ranges of timestamps
	2.14.5	Timestamp limitations
	2.14.6	Indexing
	2.14.7	Time/date components
	2.14.8	DateOffset objects
	2.14.9	Time Series-Related Instance Methods
	2.14.10	Resampling
		Time span representation
		Converting between representations
		Representing out-of-bounds spans
		Time zone handling
2.15		eltas
_,		Parsing
	2.15.2	Operations
	2.15.3	Reductions
	2.15.4	Frequency conversion
	2.15.5	Attributes
	2.15.6	TimedeltaIndex
	2.15.7	Resampling
2.16	Styling	
2.10	2.16.1	Building styles
	2.16.2	Finer control: slicing
	2.16.3	Finer Control: Display Values
	2.16.4	Builtin styles
	2.16.5	Sharing styles
	2.16.6	Other Options
	2.16.7	Fun stuff
	2.16.7	Export to Excel
	2.16.9	
2.17		Extensibility
2.17	2.17.1	· · · · · · · · · · · · · · · · · · ·
		Getting and setting options
	2.17.2	Setting startup options in Python/IPython environment
	2.17.3	Frequently Used Options
	2.17.4	Available options
	2.17.5	Number formatting
	2.17.0	· · · · · · · · · · · · · · · · · · ·
	2.17.7	Unicode formatting
2.18		
2.10		C I
	2.18.1	Cython (writing C extensions for pandas)
	2.18.2	Using Numba
2.10	2.18.3	Expression evaluation via eval()826
2.19	_	to large datasets
	2.19.1	Load less data
	2.19.2	Use efficient datatypes
	2.19.3	Use chunking
2.22	2.19.4	Use other libraries
2.20	Sparse of	lata structures

		2.20.1	SparseArray	845
		2.20.2	SparseDtype	845
		2.20.3	Sparse accessor	846
		2.20.4	Sparse calculation	846
		2.20.5		
		2.20.6		
	2.21		1 3 1	852
		2.21.1		
		2.21.2	, E	
		2.21.3	NaN, Integer NA values and NA type promotions	
		2.21.4	Differences with NumPy	
		2.21.5	Thread-safety	
		2.21.6	Byte-Ordering issues	
	2.22		ok	
	2.22	2.22.1	Idioms	
		2.22.2	Selection	
		2.22.3	MultiIndexing	
		2.22.4	Missing data	
		2.22.5	Grouping	
		2.22.6	Timeseries	
		2.22.7	8	
		2.22.8	$\epsilon$	884
		2.22.9		885
			1	890
				891
				893
		2.22.13	Creating example data	894
2	A DI -			005
3		reference		895
3	<b>API</b> 1 3.1	Input/ou	ıtput	895
3		Input/ou 3.1.1	Itput Pickling	895 895
3		Input/ou 3.1.1 3.1.2	Pickling	895 895 896
3		Input/ou 3.1.1 3.1.2 3.1.3	Pickling	895 895 896 907
3		Input/ou 3.1.1 3.1.2 3.1.3 3.1.4	Pickling	895 895 896 907 907
3		Input/ou 3.1.1 3.1.2 3.1.3 3.1.4 3.1.5	Pickling	895 895 896 907 907 913
3		Input/ou 3.1.1 3.1.2 3.1.3 3.1.4 3.1.5 3.1.6	Itput Pickling Flat file Clipboard Excel JSON HTML	895 895 896 907 907 913 919
3		Input/ou 3.1.1 3.1.2 3.1.3 3.1.4 3.1.5	ritput Pickling Flat file Clipboard Excel JSON HTML HDFStore: PyTables (HDF5)	895 895 896 907 907 913 919 921
3		Input/ou 3.1.1 3.1.2 3.1.3 3.1.4 3.1.5 3.1.6	ritput Pickling Flat file Clipboard Excel JSON HTML HDFStore: PyTables (HDF5) Feather	895 895 896 907 907 913 919 921
3		Input/ou 3.1.1 3.1.2 3.1.3 3.1.4 3.1.5 3.1.6 3.1.7	ritput Pickling Flat file Clipboard Excel JSON HTML HDFStore: PyTables (HDF5) Feather	895 895 896 907 907 913 919 921
3		Input/ou 3.1.1 3.1.2 3.1.3 3.1.4 3.1.5 3.1.6 3.1.7 3.1.8	ritput Pickling Flat file Clipboard Excel JSON HTML HDFStore: PyTables (HDF5) Feather	895 895 896 907 907 913 919 921 925
3		Input/ou 3.1.1 3.1.2 3.1.3 3.1.4 3.1.5 3.1.6 3.1.7 3.1.8 3.1.9	ritput Pickling Flat file Clipboard Excel JSON HTML HDFStore: PyTables (HDF5) Feather Parquet	895 896 907 907 913 919 921 925 926
3		Input/or 3.1.1 3.1.2 3.1.3 3.1.4 3.1.5 3.1.6 3.1.7 3.1.8 3.1.9 3.1.10	ritput Pickling Flat file Clipboard Excel JSON HTML HDFStore: PyTables (HDF5) Feather Parquet ORC	895 896 907 907 913 919 921 925 926
3		Input/or 3.1.1 3.1.2 3.1.3 3.1.4 3.1.5 3.1.6 3.1.7 3.1.8 3.1.9 3.1.10 3.1.11	ritput Pickling Flat file Clipboard Excel JSON HTML HDFStore: PyTables (HDF5) Feather Parquet ORC SAS	895 896 907 907 913 919 921 925 926 926
3		Input/ou 3.1.1 3.1.2 3.1.3 3.1.4 3.1.5 3.1.6 3.1.7 3.1.8 3.1.9 3.1.10 3.1.11 3.1.12	ritput Pickling Flat file Clipboard Excel JSON HTML HDFStore: PyTables (HDF5) Feather Parquet ORC SAS SPSS	895 895 896 907 913 919 921 925 926 927 928
3		Input/or 3.1.1 3.1.2 3.1.3 3.1.4 3.1.5 3.1.6 3.1.7 3.1.8 3.1.9 3.1.10 3.1.11 3.1.12 3.1.13	ritput Pickling Flat file Clipboard Excel JSON HTML HDFStore: PyTables (HDF5) Feather Parquet ORC SAS SPSS SQL	895 895 896 907 913 919 921 925 926 927 928 928
3		Input/ou 3.1.1 3.1.2 3.1.3 3.1.4 3.1.5 3.1.6 3.1.7 3.1.8 3.1.9 3.1.10 3.1.11 3.1.12 3.1.13 3.1.14 3.1.15	ritput Pickling Flat file Clipboard Excel JSON HTML HDFStore: PyTables (HDF5) Feather Parquet ORC SAS SPSS SPSS SQL Google BigQuery STATA	895 895 896 907 913 919 921 925 926 927 928 928 931
3	3.1	Input/ou 3.1.1 3.1.2 3.1.3 3.1.4 3.1.5 3.1.6 3.1.7 3.1.8 3.1.9 3.1.10 3.1.11 3.1.12 3.1.13 3.1.14 3.1.15	ritput Pickling Flat file Clipboard Excel JSON HTML HDFStore: PyTables (HDF5) Feather Parquet ORC SAS SPSS SQL Google BigQuery STATA functions	895 895 896 907 913 919 921 925 926 927 928 931 933
3	3.1	Input/or 3.1.1 3.1.2 3.1.3 3.1.4 3.1.5 3.1.6 3.1.7 3.1.8 3.1.9 3.1.10 3.1.11 3.1.12 3.1.13 General	ritput Pickling Flat file Clipboard Excel JSON HTML HDFStore: PyTables (HDF5) Feather Parquet ORC SAS SPSS SQL Google BigQuery STATA functions Data manipulations	895 895 896 907 913 919 921 925 926 927 928 931 933 935
3	3.1	Input/or 3.1.1 3.1.2 3.1.3 3.1.4 3.1.5 3.1.6 3.1.7 3.1.8 3.1.9 3.1.10 3.1.11 3.1.12 3.1.13 3.1.14 3.1.15 General 3.2.1 3.2.2	ritput Pickling Flat file Clipboard Excel JSON HTML HDFStore: PyTables (HDF5) Feather Parquet ORC SAS SPSS SQL Google BigQuery STATA functions Data manipulations Top-level missing data	895 895 896 907 913 919 921 925 926 927 928 933 933 935 965
3	3.1	Input/or 3.1.1 3.1.2 3.1.3 3.1.4 3.1.5 3.1.6 3.1.7 3.1.8 3.1.9 3.1.10 3.1.11 3.1.12 3.1.13 3.1.14 3.1.15 General 3.2.1 3.2.2 3.2.3	atput Pickling Flat file Clipboard Excel JSON HTML HDFStore: PyTables (HDF5) Feather Parquet ORC SAS SPSS SQL Google BigQuery STATA functions Data manipulations Top-level missing data Top-level conversions	895 895 896 907 913 919 921 925 926 927 928 931 933 935 935 971
3	3.1	Input/or 3.1.1 3.1.2 3.1.3 3.1.4 3.1.5 3.1.6 3.1.7 3.1.8 3.1.9 3.1.10 3.1.11 3.1.12 3.1.13 3.1.14 3.1.15 General 3.2.1 3.2.2 3.2.3 3.2.4	rtput Pickling Flat file Clipboard Excel JSON HTML HDFStore: PyTables (HDF5) Feather Parquet ORC SAS SPSS SQL Google BigQuery STATA functions Data manipulations Top-level missing data Top-level conversions Top-level dealing with datetimelike	895 895 896 907 913 919 921 925 926 927 928 931 933 935 935 971 973
3	3.1	Input/or 3.1.1 3.1.2 3.1.3 3.1.4 3.1.5 3.1.6 3.1.7 3.1.8 3.1.9 3.1.10 3.1.11 3.1.12 3.1.13 3.1.14 3.1.15 General 3.2.1 3.2.2 3.2.3 3.2.4 3.2.5	rtput Pickling Flat file Clipboard Excel JSON HTML HDFStore: PyTables (HDF5) Feather Parquet ORC SAS SPSS SQL Google BigQuery STATA functions Data manipulations Top-level missing data Top-level dealing with datetimelike Top-level dealing with intervals	895 895 896 907 913 919 921 925 926 927 928 931 933 935 965 971 973 982
3	3.1	Input/or 3.1.1 3.1.2 3.1.3 3.1.4 3.1.5 3.1.6 3.1.7 3.1.8 3.1.9 3.1.10 3.1.11 3.1.12 3.1.13 3.1.14 3.1.15 General 3.2.1 3.2.2 3.2.3 3.2.4 3.2.5 3.2.6	rtput Pickling Flat file Clipboard Excel JSON HTML HDFStore: PyTables (HDF5) Feather Parquet ORC SAS SPSS SQL Google BigQuery STATA functions Data manipulations Top-level missing data Top-level dealing with datetimelike Top-level dealing with intervals Top-level evaluation	895 895 896 907 913 919 921 925 926 927 928 931 933 935 965 971 973 982 984
3	3.1	Input/or 3.1.1 3.1.2 3.1.3 3.1.4 3.1.5 3.1.6 3.1.7 3.1.8 3.1.9 3.1.10 3.1.11 3.1.12 3.1.13 3.1.14 3.1.15 General 3.2.1 3.2.2 3.2.3 3.2.4 3.2.5	riput Pickling Flat file Clipboard Excel JSON HTML HDFStore: PyTables (HDF5) Feather Parquet ORC SAS SPSS SQL Google BigQuery STATA functions Data manipulations Top-level missing data Top-level dealing with datetimelike Top-level evaluation Hashing	895 895 896 907 913 919 921 925 926 927 928 931 933 935 965 971 973 982

3.3	Series .	987
	3.3.1	Constructor
	3.3.2	Attributes
	3.3.3	Conversion
	3.3.4	Indexing, iteration
	3.3.5	Binary operator functions
	3.3.6	Function application, groupby & window
	3.3.7	Computations / descriptive stats
	3.3.8	Reindexing / selection / label manipulation
	3.3.9	Missing data handling
	3.3.10	Reshaping, sorting
	3.3.11	Combining / joining / merging
	3.3.12	Time series-related
	3.3.13	Accessors
	3.3.14	Plotting
2.4	3.3.15	Serialization / IO / conversion
3.4		nme
	3.4.1	Constructor
	3.4.2	Attributes and underlying data
	3.4.3 3.4.4	Conversion
	3.4.4	Indexing, iterationBinary operator functions
	3.4.6	Function application, GroupBy & window
	3.4.7	Computations / descriptive stats
	3.4.7	Reindexing / selection / label manipulation
	3.4.9	Missing data handling
	3.4.10	Reshaping, sorting, transposing
	3.4.11	Combining / joining / merging
	3.4.12	Time series-related
	3.4.13	Metadata
	3.4.14	Plotting
	3.4.15	Sparse accessor
	3.4.16	Serialization / IO / conversion
3.5		arrays
	3.5.1	pandas.array
	3.5.2	Datetime data
	3.5.3	Timedelta data
	3.5.4	Timespan data
	3.5.5	Period
	3.5.6	Interval data
	3.5.7	Nullable integer
	3.5.8	Categorical data
	3.5.9	Sparse data
	3.5.10	Text data
	3.5.11	Boolean data with missing values
3.6		1810
3.7		bjects
	3.7.1	Index
	3.7.2	Numeric Index
	3.7.3	CategoricalIndex
	3.7.4	IntervalIndex
	3.7.5	MultiIndex
	3.7.6	DatetimeIndex
	3.7.7	TimedeltaIndex

	3.7.8	PeriodIndex
3.8		Seets
0.0	3.8.1	DateOffset
	3.8.2	BusinessDay
	3.8.3	BusinessHour
	3.8.4	CustomBusinessDay
	3.8.5	CustomBusinessHour
	3.8.6	
		MonthOffset         1971           MonthEnd         1975
	3.8.7	
	3.8.8	MonthBegin
	3.8.9	BusinessMonthEnd
	3.8.10	BusinessMonthBegin
	3.8.11	CustomBusinessMonthEnd
	3.8.12	CustomBusinessMonthBegin
	3.8.13	SemiMonthOffset
	3.8.14	SemiMonthEnd
	3.8.15	SemiMonthBegin
	3.8.16	Week
	3.8.17	WeekOfMonth
	3.8.18	LastWeekOfMonth
	3.8.19	QuarterOffset
	3.8.20	BQuarterEnd
	3.8.21	BQuarterBegin
	3.8.22	QuarterEnd
	3.8.23	QuarterBegin
	3.8.24	YearOffset
	3.8.25	BYearEnd
	3.8.26	BYearBegin
	3.8.27	YearEnd
	3.8.28	YearBegin
	3.8.29	FY5253
	3.8.30	FY5253Quarter
	3.8.31	Easter
	3.8.32	Tick
	3.8.33	Day
	3.8.34	Hour
	3.8.35	Minute
	3.8.36	Second
	3.8.37	Milli
	3.8.38	Micro
	3.8.39	Nano
	3.8.40	BDay
	3.8.41	BMonthEnd
	3.8.42	BMonthBegin
	3.8.43	CBMonthEnd
	3.8.44	
2.0	3.8.45	CDay
3.9		ncies
0.10	3.9.1	pandas.tseries.frequencies.to_offset
3.10		7
	3.10.1	Standard moving window functions
	3.10.2	Standard expanding window functions
	3.10.3	Exponentially-weighted moving window functions
	3.10.4	Window Indexer

3.11		By
		Indexing, iteration
	3.11.2	Function application
	3.11.3	Computations / descriptive stats
3.12	Resamp	oling
	3.12.1	Indexing, iteration
	3.12.2	Function application
	3.12.3	Upsampling
	3.12.4	Computations / descriptive stats
3.13	Style	
	3.13.1	Styler constructor
	3.13.2	Styler properties
	3.13.3	Style application
	3.13.4	Builtin styles
	3.13.5	Style export and import
3.14		g
011.	3.14.1	pandas.plotting.andrews_curves
	3.14.2	pandas.plotting.autocorrelation_plot
	3.14.3	pandas.plotting.bootstrap_plot
	3.14.4	pandas.plotting.boxplot
	3.14.5	pandas.plotting.deregister_matplotlib_converters
	3.14.6	pandas.plotting.lag_plot
	3.14.7	pandas.plotting.parallel_coordinates
	3.14.7	pandas.plotting.plot_params
	3.14.9	pandas.plotting.radviz
		pandas.plotting.register_matplotlib_converters
		pandas.plotting.scatter_matrix
2 15		pandas.plotting.table
3.15		l utility functions
	3.15.1	Working with options
	3.15.2	Testing functions
	3.15.3	Exceptions and warnings
	3.15.4	Data types related functionality
3.16		ons
	3.16.1	pandas.api.extensions.register_extension_dtype
	3.16.2	pandas.api.extensions.register_dataframe_accessor
	3.16.3	pandas.api.extensions.register_series_accessor
		pandas.api.extensions.register_index_accessor
	3.16.5	pandas.api.extensions.ExtensionDtype
	3.16.6	pandas.api.extensions.ExtensionArray
	3.16.7	pandas.arrays.PandasArray
	3.16.8	pandas.api.indexers.check_array_indexer
		***
	lopment	
4.1		outing to pandas
	4.1.1	Where to start?
	4.1.2	Bug reports and enhancement requests
	4.1.3	Working with the code
	4.1.4	Contributing to the documentation
	4.1.5	Contributing to the code base
	4.1.6	Contributing your changes to <i>pandas</i>
4.2		code style guide
	4.2.1	Patterns
	4.2.2	String formatting

4

4.3	Pandas	Maintenance													
	4.3.1	Roles													
	4.3.2	Tasks													
	4.3.3	Issue Triage													
	4.3.4	Closing Issues													
	4.3.5	Reviewing Pull Requests													
	4.3.6	Cleaning up old Issues													
	4.3.7	Cleaning up old Pull Requests													
	4.3.8	Becoming a pandas maintainer													
4.4	Internal														
	4.4.1	Indexing													
	4.4.2	Subclassing pandas data structures													
4.5		ng pandas													
	4.5.1	Registering custom accessors													
	4.5.2	Extension types													
	4.5.3	Subclassing pandas data structures													
	4.5.4	Plotting backends													
4.6	-	per													
4.7	4.6.1	Storing pandas DataFrame objects in Apa		-											
4.7		Waster Delta													
	4.7.1	Version Policy													
4.0	4.7.2	Python Support													
4.8	4.8.1	ap													
	4.8.2	String data type													
	4.8.3	Apache Arrow interoperability													
	4.8.4	Block manager rewrite													
	4.8.5	Decoupling of indexing and internals													
	4.8.6	Numba-accelerated operations													
	4.8.7	Documentation improvements													
	4.8.8	Package docstring validation													
	4.8.9	Performance monitoring													
	4.8.10	Roadmap Evolution													
4.9		per Meetings													
7.7	4.9.1	Minutes													
	4.9.2	Calendar													
		Culcitati	 		 • •	• • •	• •	• •	• •	• •	•	•	• •	•	_500
Relea	ase Notes	5												23	389
5.1	Version	1.0	 		 									. 2	2389
	5.1.1	What's new in 1.0.5 (June 17, 2020)	 		 									. 2	2389
	5.1.2	What's new in 1.0.4 (May 28, 2020)	 		 									. 2	2390
	5.1.3	What's new in 1.0.3 (March 17, 2020)	 		 									. 2	2391
	5.1.4	What's new in 1.0.2 (March 12, 2020)	 		 									. 2	2392
	5.1.5	What's new in 1.0.1 (February 5, 2020)	 		 									. 2	2395
	5.1.6	What's new in 1.0.0 (January 29, 2020)	 		 									. 2	2398
5.2	Version														2437
	5.2.1	What's new in 0.25.3 (October 31, 2019)													2437
	5.2.2	What's new in 0.25.2 (October 15, 2019)													2437
	5.2.3	What's new in 0.25.1 (August 21, 2019)													2438
	5.2.4	What's new in 0.25.0 (July 18, 2019)													2441
5.3	Version														2479
	5.3.1	Whats new in 0.24.2 (March 12, 2019)													2479
	5.3.2	Whats new in 0.24.1 (February 3, 2019)													2482
	5.3.3	What's new in 0.24.0 (January 25, 2019)	 		 									. 2	2483

5

5.4	Version	0.23
	5.4.1	What's new in 0.23.4 (August 3, 2018)
	5.4.2	What's new in 0.23.3 (July 7, 2018)
	5.4.3	What's new in 0.23.2 (July 5, 2018)
	5.4.4	What's new in 0.23.1 (June 12, 2018)
	5.4.5	What's new in 0.23.0 (May 15, 2018)
5.5	Version	
3.3	5.5.1	v0.22.0 (December 29, 2017)
5.6	Version	
5.0	5.6.1	v0.21.1 (December 12, 2017)
	5.6.2	v0.21.1 (December 12, 2017)
5.7		0.20
5.7		
	5.7.1	v0.20.3 (July 7, 2017)
	5.7.2	v0.20.2 (June 4, 2017)
<b>~</b> 0	5.7.3	v0.20.1 (May 5, 2017)
5.8		0.19
	5.8.1	v0.19.2 (December 24, 2016)
	5.8.2	v0.19.1 (November 3, 2016)
	5.8.3	v0.19.0 (October 2, 2016)
5.9	Version	
	5.9.1	v0.18.1 (May 3, 2016)
	5.9.2	v0.18.0 (March 13, 2016)
5.10	Version	0.17
	5.10.1	v0.17.1 (November 21, 2015)
	5.10.2	v0.17.0 (October 9, 2015)
5.11	Version	0.16
	5.11.1	v0.16.2 (June 12, 2015)
	5.11.2	v0.16.1 (May 11, 2015)
	5.11.3	v0.16.0 (March 22, 2015)
5.12	Version	
	5.12.1	v0.15.2 (December 12, 2014)
	5.12.2	v0.15.1 (November 9, 2014)
	5.12.3	v0.15.0 (October 18, 2014)
5.13	Version	
5.15	5.13.1	v0.14.1 (July 11, 2014)
	5.13.2	v0.14.0 (May 31, 2014)
5 14		0.13
5.14		v0.13.1 (February 3, 2014)
	5.14.1	v0.13.1 (February 3, 2014)
5.15	Version	
3.13		
F 16	5.15.1	v0.12.0 (July 24, 2013)
5.16	Version	
5 15	5.16.1	v0.11.0 (April 22, 2013)
5.17		0.10
	5.17.1	v0.10.1 (January 22, 2013)
	5.17.2	v0.10.0 (December 17, 2012)
5.18		0.9
	5.18.1	v0.9.1 (November 14, 2012)
	5.18.2	v0.9.0 (October 7, 2012)
5.19	Version	0.8
	5.19.1	v0.8.1 (July 22, 2012)
	5.19.2	v0.8.0 (June 29, 2012)
5.20		0.7
	5.20.1	v.0.7.3 (April 12, 2012)

Python N	Module Index 36	073				
Bibliography 3						
	5.23.1 v.0.4.1 through v0.4.3 (September 25 - October 9, 2011)	3068				
5.23	Version 0.4	3068				
	5.22.1 v.0.5.0 (October 24, 2011)	3066				
5.22	Version 0.5	3066				
	5.21.2 v.0.6.0 (November 25, 2011)	3064				
	5.21.1 v.0.6.1 (December 13, 2011)	3063				
5.21	Version 0.6	3063				
	5.20.4 v.0.7.0 (February 9, 2012)	3057				
	5.20.3 v.0.7.1 (February 29, 2012)	3056				
	5.20.2 v.0.7.2 (March 16, 2012)	3055				

**Date**: Jun 17, 2020 **Version**: 1.0.5

**Download documentation**: PDF Version | Zipped HTML

Useful links: Binary Installers | Source Repository | Issues & Ideas | Q&A Support | Mailing List

pandas is an open source, BSD-licensed library providing high-performance, easy-to-use data structures and data analysis tools for the Python programming language.

To the getting started guides

To the user guide

To the reference guide

To the development guide

CONTENTS 1

2 CONTENTS

**CHAPTER** 

**ONE** 

## **GETTING STARTED**

### 1.1 Installation

Before you can use pandas, you'll need to get it installed.

Pandas is part of the Anaconda distribution and can be installed with Anaconda or Miniconda:

conda install pandas

Pandas can be installed via pip from PyPI.

pip install pandas

Learn more

## 1.2 Intro to pandas

Straight to tutorial...

When working with tabular data, such as data stored in spreadsheets or databases, Pandas is the right tool for you. Pandas will help you to explore, clean and process your data. In Pandas, a data table is called a <code>DataFrame</code>.

To introduction tutorial

To user guide

Straight to tutorial...

Pandas supports the integration with many file formats or data sources out of the box (csv, excel, sql, json, parquet,...). Importing data from each of these data sources is provided by function with the prefix read\_\*. Similarly, the to\_\* methods are used to store data.

To introduction tutorial

To user guide

Straight to tutorial...

Selecting or filtering specific rows and/or columns? Filtering the data on a condition? Methods for slicing, selecting, and extracting the data you need are available in Pandas.

To introduction tutorial

To user guide

Straight to tutorial...

Pandas provides plotting your data out of the box, using the power of Matplotlib. You can pick the plot type (scatter, bar, boxplot,...) corresponding to your data.

To introduction tutorial

To user guide

Straight to tutorial...

There is no need to loop over all rows of your data table to do calculations. Data manipulations on a column work elementwise. Adding a column to a <code>DataFrame</code> based on existing data in other columns is straightforward.

To introduction tutorial

To user guide

Straight to tutorial...

Basic statistics (mean, median, min, max, counts...) are easily calculable. These or custom aggregations can be applied on the entire data set, a sliding window of the data or grouped by categories. The latter is also known as the split-apply-combine approach.

To introduction tutorial

To user guide

Straight to tutorial...

Change the structure of your data table in multiple ways. You can melt () your data table from wide to long/tidy form or pivot () from long to wide format. With aggregations built-in, a pivot table is created with a single command.

To introduction tutorial

To user guide

Straight to tutorial...

Multiple tables can be concatenated both column wise as row wise and database-like join/merge operations are provided to combine multiple tables of data.

To introduction tutorial

To user guide

Straight to tutorial...

Pandas has great support for time series and has an extensive set of tools for working with dates, times, and time-indexed data.

To introduction tutorial

To user guide

Straight to tutorial...

Data sets do not only contain numerical data. Pandas provides a wide range of functions to cleaning textual data and extract useful information from it.

To introduction tutorial

To user guide

## 1.3 Coming from...

Currently working with other software for data manipulation in a tabular format? You're probably familiar to typical data operations and know *what* to do with your tabular data, but lacking the syntax to execute these operations. Get to know the pandas syntax by looking for equivalents from the software you already know:

Learn more

Learn more

Learn more

Learn more

## 1.4 Community tutorials

The community produces a wide variety of tutorials available online. Some of the material is enlisted in the community contributed *Tutorials*.

#### 1.4.1 Installation

The easiest way to install pandas is to install it as part of the Anaconda distribution, a cross platform distribution for data analysis and scientific computing. This is the recommended installation method for most users.

Instructions for installing from source, PyPI, ActivePython, various Linux distributions, or a development version are also provided.

#### Python version support

Officially Python 3.6.1 and above, 3.7, and 3.8.

#### Installing pandas

#### Installing with Anaconda

Installing pandas and the rest of the NumPy and SciPy stack can be a little difficult for inexperienced users.

The simplest way to install not only pandas, but Python and the most popular packages that make up the SciPy stack (IPython, NumPy, Matplotlib, ...) is with Anaconda, a cross-platform (Linux, Mac OS X, Windows) Python distribution for data analytics and scientific computing.

After running the installer, the user will have access to pandas and the rest of the SciPy stack without needing to install anything else, and without needing to wait for any software to be compiled.

Installation instructions for Anaconda can be found here.

A full list of the packages available as part of the Anaconda distribution can be found here.

Another advantage to installing Anaconda is that you don't need admin rights to install it. Anaconda can install in the user's home directory, which makes it trivial to delete Anaconda if you decide (just delete that folder).

#### **Installing with Miniconda**

The previous section outlined how to get pandas installed as part of the Anaconda distribution. However this approach means you will install well over one hundred packages and involves downloading the installer which is a few hundred megabytes in size.

If you want to have more control on which packages, or have a limited internet bandwidth, then installing pandas with Miniconda may be a better solution.

Conda is the package manager that the Anaconda distribution is built upon. It is a package manager that is both cross-platform and language agnostic (it can play a similar role to a pip and virtualenv combination).

Miniconda allows you to create a minimal self contained Python installation, and then use the Conda command to install additional packages.

First you will need Conda to be installed and downloading and running the Miniconda will do this for you. The installer can be found here

The next step is to create a new conda environment. A conda environment is like a virtualenv that allows you to specify a specific version of Python and set of libraries. Run the following commands from a terminal window:

```
conda create -n name_of_my_env python
```

This will create a minimal environment with only Python installed in it. To put your self inside this environment run:

```
source activate name_of_my_env
```

On Windows the command is:

```
activate name_of_my_env
```

The final step required is to install pandas. This can be done with the following command:

```
conda install pandas
```

To install a specific pandas version:

```
conda install pandas=0.20.3
```

To install other packages, IPython for example:

```
conda install ipython
```

To install the full Anaconda distribution:

```
conda install anaconda
```

If you need packages that are available to pip but not conda, then install pip, and then use pip to install those packages:

```
conda install pip
pip install django
```

#### Installing from PyPI

pandas can be installed via pip from PyPI.

```
pip install pandas
```

#### Installing with ActivePython

Installation instructions for ActivePython can be found here. Versions 2.7, 3.5 and 3.6 include pandas.

#### Installing using your Linux distribution's package manager.

The commands in this table will install pandas for Python 3 from your distribution. To install pandas for Python 2, you may need to use the python-pandas package.

Distribution Status		Download / Reposi-	Install method
		tory Link	
Debian	stable	official Debian reposi-	sudo apt-get install python3-pandas
		tory	
Debian &	unstable	NeuroDebian	sudo apt-get install python3-pandas
Ubuntu	(latest		
	pack-		
	ages)		
Ubuntu	stable	official Ubuntu reposi-	sudo apt-get install python3-pandas
		tory	
OpenSuse	stable	OpenSuse Repository	zypper in python3-pandas
Fedora	stable	official Fedora reposi-	dnf install python3-pandas
		tory	
Centos/RH	E <b>k</b> table	EPEL repository	yum install python3-pandas

**However**, the packages in the linux package managers are often a few versions behind, so to get the newest version of pandas, it's recommended to install using the pip or conda methods described above.

#### Installing from source

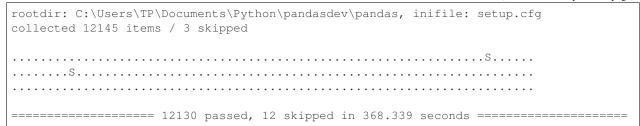
See the *contributing guide* for complete instructions on building from the git source tree. Further, see *creating a development environment* if you wish to create a *pandas* development environment.

#### Running the test suite

pandas is equipped with an exhaustive set of unit tests, covering about 97% of the code base as of this writing. To run it on your machine to verify that everything is working (and that you have all of the dependencies, soft and hard, installed), make sure you have pytest >= 5.0.1 and Hypothesis >= 3.58, then run:

(continues on next page)

(continued from previous page)



#### **Dependencies**

Package	Minimum supported version
setuptools	24.2.0
NumPy	1.13.3
python-dateutil	2.6.1
pytz	2017.2

#### **Recommended dependencies**

- numexpr: for accelerating certain numerical operations. numexpr uses multiple cores as well as smart chunking and caching to achieve large speedups. If installed, must be Version 2.6.2 or higher.
- bottleneck: for accelerating certain types of nan evaluations. bottleneck uses specialized cython routines to achieve large speedups. If installed, must be Version 1.2.1 or higher.

**Note:** You are highly encouraged to install these libraries, as they provide speed improvements, especially when working with large data sets.

#### **Optional dependencies**

Pandas has many optional dependencies that are only used for specific methods. For example, <code>pandas.read\_hdf()</code> requires the <code>pytables</code> package, while <code>DataFrame.to\_markdown()</code> requires the tabulate package. If the optional dependency is not installed, pandas will raise an <code>ImportError</code> when the method requiring that dependency is called.

Dependency	Minimum Version	Notes
BeautifulSoup4	4.6.0	HTML parser for read_html (see <i>note</i> )
Jinja2		Conditional formatting with DataFrame.style
PyQt4		Clipboard I/O
PyQt5		Clipboard I/O
PyTables	3.4.2	HDF5-based reading / writing
SQLAlchemy	1.1.4	SQL support for databases other than sqlite
SciPy	0.19.0	Miscellaneous statistical functions
XLsxWriter	0.9.8	Excel writing
blosc		Compression for HDF5
fastparquet	0.3.2	Parquet reading / writing

continues on next page

Table 1 – continued from previous page

Dependency	Minimum Version	Notes
gcsfs	0.2.2	Google Cloud Storage access
html5lib		HTML parser for read_html (see <i>note</i> )
lxml	3.8.0	HTML parser for read_html (see <i>note</i> )
matplotlib	2.2.2	Visualization
numba	0.46.0	Alternative execution engine for rolling operations
openpyxl	2.5.7	Reading / writing for xlsx files
pandas-gbq	0.8.0	Google Big Query access
psycopg2		PostgreSQL engine for sqlalchemy
pyarrow	0.12.0	Parquet, ORC (requires 0.13.0), and feather reading / writing
pymysql	0.7.11	MySQL engine for sqlalchemy
pyreadstat		SPSS files (.sav) reading
pytables	3.4.2	HDF5 reading / writing
pyxlsb	1.0.6	Reading for xlsb files
qtpy		Clipboard I/O
s3fs	0.3.0	Amazon S3 access
tabulate	0.8.3	Printing in Markdown-friendly format (see tabulate)
xarray	0.8.2	pandas-like API for N-dimensional data
xclip		Clipboard I/O on linux
xlrd	1.1.0	Excel reading
xlwt	1.2.0	Excel writing
xsel		Clipboard I/O on linux
zlib		Compression for HDF5

#### Optional dependencies for parsing HTML

One of the following combinations of libraries is needed to use the top-level read\_html() function:

Changed in version 0.23.0.

- BeautifulSoup4 and html5lib
- BeautifulSoup4 and lxml
- BeautifulSoup4 and html5lib and lxml
- Only lxml, although see HTML Table Parsing for reasons as to why you should probably **not** take this approach.

#### Warning:

- if you install BeautifulSoup4 you must install either lxml or html5lib or both. read\_html() will **not** work with *only* BeautifulSoup4 installed.
- You are highly encouraged to read *HTML Table Parsing gotchas*. It explains issues surrounding the installation and usage of the above three libraries.

#### 1.4.2 Package overview

pandas is a Python package providing fast, flexible, and expressive data structures designed to make working with "relational" or "labeled" data both easy and intuitive. It aims to be the fundamental high-level building block for doing practical, real world data analysis in Python. Additionally, it has the broader goal of becoming the most powerful and flexible open source data analysis / manipulation tool available in any language. It is already well on its way toward this goal.

pandas is well suited for many different kinds of data:

- Tabular data with heterogeneously-typed columns, as in an SQL table or Excel spreadsheet
- Ordered and unordered (not necessarily fixed-frequency) time series data.
- · Arbitrary matrix data (homogeneously typed or heterogeneous) with row and column labels
- Any other form of observational / statistical data sets. The data actually need not be labeled at all to be placed into a pandas data structure

The two primary data structures of pandas, <code>Series</code> (1-dimensional) and <code>DataFrame</code> (2-dimensional), handle the vast majority of typical use cases in finance, statistics, social science, and many areas of engineering. For R users, <code>DataFrame</code> provides everything that R's data.frame provides and much more. pandas is built on top of NumPy and is intended to integrate well within a scientific computing environment with many other 3rd party libraries.

Here are just a few of the things that pandas does well:

- Easy handling of **missing data** (represented as NaN) in floating point as well as non-floating point data
- Size mutability: columns can be inserted and deleted from DataFrame and higher dimensional objects
- Automatic and explicit **data alignment**: objects can be explicitly aligned to a set of labels, or the user can simply ignore the labels and let *Series*, *DataFrame*, etc. automatically align the data for you in computations
- Powerful, flexible **group by** functionality to perform split-apply-combine operations on data sets, for both aggregating and transforming data
- Make it **easy to convert** ragged, differently-indexed data in other Python and NumPy data structures into DataFrame objects
- Intelligent label-based slicing, fancy indexing, and subsetting of large data sets
- Intuitive **merging** and **joining** data sets
- Flexible reshaping and pivoting of data sets
- Hierarchical labeling of axes (possible to have multiple labels per tick)
- Robust IO tools for loading data from flat files (CSV and delimited), Excel files, databases, and saving / loading data from the ultrafast HDF5 format
- **Time series**-specific functionality: date range generation and frequency conversion, moving window statistics, date shifting and lagging.

Many of these principles are here to address the shortcomings frequently experienced using other languages / scientific research environments. For data scientists, working with data is typically divided into multiple stages: munging and cleaning data, analyzing / modeling it, then organizing the results of the analysis into a form suitable for plotting or tabular display, pandas is the ideal tool for all of these tasks.

Some other notes

• pandas is **fast**. Many of the low-level algorithmic bits have been extensively tweaked in Cython code. However, as with anything else generalization usually sacrifices performance. So if you focus on one feature for your application you may be able to create a faster specialized tool.

- pandas is a dependency of statsmodels, making it an important part of the statistical computing ecosystem in Python.
- pandas has been used extensively in production in financial applications.

#### **Data structures**

Dimensions	Name	Description
1	Series	1D labeled homogeneously-typed array
2	DataFrame	General 2D labeled, size-mutable tabular structure with potentially
		heterogeneously-typed column

#### Why more than one data structure?

The best way to think about the pandas data structures is as flexible containers for lower dimensional data. For example, DataFrame is a container for Series, and Series is a container for scalars. We would like to be able to insert and remove objects from these containers in a dictionary-like fashion.

Also, we would like sensible default behaviors for the common API functions which take into account the typical orientation of time series and cross-sectional data sets. When using ndarrays to store 2- and 3-dimensional data, a burden is placed on the user to consider the orientation of the data set when writing functions; axes are considered more or less equivalent (except when C- or Fortran-contiguousness matters for performance). In pandas, the axes are intended to lend more semantic meaning to the data; i.e., for a particular data set there is likely to be a "right" way to orient the data. The goal, then, is to reduce the amount of mental effort required to code up data transformations in downstream functions.

For example, with tabular data (DataFrame) it is more semantically helpful to think of the **index** (the rows) and the **columns** rather than axis 0 and axis 1. Iterating through the columns of the DataFrame thus results in more readable code:

```
for col in df.columns:
    series = df[col]
    # do something with series
```

#### Mutability and copying of data

All pandas data structures are value-mutable (the values they contain can be altered) but not always size-mutable. The length of a Series cannot be changed, but, for example, columns can be inserted into a DataFrame. However, the vast majority of methods produce new objects and leave the input data untouched. In general we like to **favor immutability** where sensible.

#### **Getting support**

The first stop for pandas issues and ideas is the Github Issue Tracker. If you have a general question, pandas community experts can answer through Stack Overflow.