

---

# **pygdf Documentation**

*Release 0.1.0*

**Continuum Analytics**

**Jan 23, 2018**



---

## Contents:

---

<b>1</b>	<b>API Reference</b>	<b>1</b>
1.1	DataFrame . . . . .	1
1.2	Series . . . . .	9
1.3	Groupby . . . . .	14
1.4	GpuArrowReader . . . . .	15
<b>2</b>	<b>Indices and tables</b>	<b>17</b>



## 1.1 DataFrame

**class** `pygdf.dataframe.DataFrame` (*name\_series=None*)  
A GPU Dataframe object.

### Examples

Build dataframe with `__setitem__`

```
>>> from pygdf.dataframe import DataFrame
>>> df = DataFrame()
>>> df['key'] = [0, 1, 2, 3, 4]
>>> df['val'] = [float(i + 10) for i in range(5)] # insert column
>>> df
  key val
0 0  10.0
1 1  11.0
2 2  12.0
3 3  13.0
4 4  14.0
>>> len(df)
5
```

Build dataframe with initializer

```
>>> import numpy as np
>>> df2 = DataFrame([('a', np.arange(10)),
...                 ('b', np.random.random(10))])
>>> df2
  a b
0 0 0.777831724018
1 1 0.604480034669
```

```
2 2 0.664111858618
3 3 0.887777513028
4 4 0.55838311246
[5 more rows]
```

Convert from a Pandas DataFrame.

```
>>> import pandas as pd
>>> from pygdf.dataframe import DataFrame
>>> pdf = pd.DataFrame({'a': [0, 1, 2, 3],
...                    'b': [0.1, 0.2, None, 0.3]})
>>> pdf
a    b
0  0  0.1
1  1  0.2
2  2  NaN
3  3  0.3
>>> df = DataFrame.from_pandas(pdf)
>>> df
a  b
0  0  0.1
1  1  0.2
2  2  nan
3  3  0.3
```

## Attributes

<i>columns</i>	Returns a tuple of columns
<i>dtypes</i>	Return the dtypes in this object.
<i>index</i>	Returns the index of the DataFrame
<i>loc</i>	Returns a label-based indexer for row-slicing and column selection.

## Methods

<i>add_column(name, data)</i>	Add a column
<i>apply_chunks(func, incol, outcols[, ...])</i>	Transform user-specified chunks using the user-provided function.
<i>apply_rows(func, incol, outcols, kwargs)</i>	Transform each row using the user-provided function.
<i>as_gpu_matrix([columns])</i>	Convert to a matrix in device memory.
<i>as_matrix([columns])</i>	Covert to a matrix in host memory.
<i>copy()</i>	Shallow copy this dataframe
<i>drop_column(name)</i>	Drop a column by <i>name</i>
<i>from_pandas(dataframe)</i>	Convert from a Pandas DataFrame.
<i>from_records(data[, index, columns])</i>	Convert from a numpy recarray or structured array.
<i>groupby(by[, sort, as_index])</i>	Groupby
<i>head([n])</i>	
<i>join(other[, on, how, lsuffix, rsuffix, sort])</i>	Join columns with other DataFrame on index or on a key column.
<i>label_encoding(column, prefix, cats[, ...])</i>	Encode labels in a column with label encoding.

Continued on next page

Table 1.2 – continued from previous page

<code>nlargest(n, columns[, keep])</code>	Get the rows of the DataFrame sorted by the <code>n</code> largest value of <code>columns</code>
<code>nsmallest(n, columns[, keep])</code>	Get the rows of the DataFrame sorted by the <code>n</code> smallest value of <code>columns</code>
<code>one_hot_encoding(column, prefix, cats[, ...])</code>	Expand a column with one-hot-encoding.
<code>query(expr)</code>	Query with a boolean expression using Numba to compile a GPU kernel.
<code>reset_index()</code>	
<code>set_index(index)</code>	Return a new DataFrame with a new index
<code>sort_index([ascending])</code>	Sort by the index
<code>sort_values(by[, ascending])</code>	Sort by values.
<code>take(positions[, ignore_index])</code>	
<code>to_pandas()</code>	Convert to a Pandas DataFrame.
<code>to_records([index])</code>	Covert to a numpy recarray
<code>to_string([nrows, ncols])</code>	Convert to string

**add\_column** (*name, data*)

Add a column

**Parameters name** : str

Name of column to be added.

**data** : Series, array-like

Values to be added.

**apply\_chunks** (*func, incol, outcol, kwargs={}, chunks=None, tpb=1*)

Transform user-specified chunks using the user-provided function.

**Parameters func** : function

The transformation function that will be executed on the CUDA GPU.

**incol**: list

A list of names of input columns.

**outcol**: dict

A dictionary of output column names and their dtype.

**kwargs**: dict

name-value of extra arguments. These values are passed directly into the function.

**chunks** : int or Series-like

If it is an `int`, it is the chunksize. If it is an array, it contains integer offset for the start of each chunk. The span of a chunk for chunk `i`-th is `data[chunks[i] : chunks[i + 1]]` for any `i + 1 < chunks.size`; or, `data[chunks[i] : ]` for the `i == len(chunks) - 1`.

**tpb** : int; optional

It is the thread-per-block for the underlying kernel. The default uses 1 thread to emulate serial execution for each chunk. It is a good starting point but inefficient. Its maximum possible value is limited by the available CUDA GPU resources.

**See also:**

`apply_rows`

## Examples

For `tpb > 1`, `func` is executed by `tpb` number of threads concurrently. To access the thread id and count, use `numba.cuda.threadIdx.x` and `numba.cuda.blockDim.x`, respectively (See [numba CUDA kernel documentation](#)).

In the example below, the *kernel* is invoked concurrently on each specified chunk. The *kernel* computes the corresponding output for the chunk. By looping over the range `range(cuda.threadIdx.x, in1.size, cuda.blockDim.x)`, the *kernel* function can be used with any *tpb* in a efficient manner.

```
>>> from numba import cuda
>>> def kernel(in1, in2, in3, out1):
...     for i in range(cuda.threadIdx.x, in1.size, cuda.blockDim.x):
...         x = in1[i]
...         y = in2[i]
...         z = in3[i]
...         out1[i] = x * y + z
```

**apply\_rows** (*func, incol, outcols, kwargs*)

Transform each row using the user-provided function.

**Parameters** `func` : function

The transformation function that will be executed on the CUDA GPU.

**incols**: list

A list of names of input columns.

**outcols**: dict

A dictionary of output column names and their dtype.

**kwargs**: dict

name-value of extra arguments. These values are passed directly into the function.

## Examples

With a DataFrame like so:

```
>>> df = DataFrame()
>>> df['in1'] = in1 = np.arange(nelem)
>>> df['in2'] = in2 = np.arange(nelem)
>>> df['in3'] = in3 = np.arange(nelem)
```

Define the user function for `.apply_rows`:

```
>>> def kernel(in1, in2, in3, out1, out2, extra1, extra2):
...     for i, (x, y, z) in enumerate(zip(in1, in2, in3)):
...         out1[i] = extra2 * x - extra1 * y
...         out2[i] = y - extra1 * z
```

The user function should loop over the columns and set the output for each row. Each iteration of the loop **MUST** be independent of each other. The order of the loop execution can be arbitrary.

Call `.apply_rows` with the name of the input columns, the name and dtype of the output columns, and, optionally, a dict of extra arguments.



```

>>> outdf = df.apply_rows(kernel,
...                        incolcols=['in1', 'in2', 'in3'],
...                        outcols=dict(out1=np.float64,
...                                     out2=np.float64),
...                        kwargs=dict(extra1=2.3, extra2=3.4))

```

**Notes**

When `func` is invoked, the array args corresponding to the input/output are strided in a way that improves parallelism on the GPU. The loop in the function may look like serial code but it will be executed concurrently by multiple threads.

**as\_gpu\_matrix** (*columns=None*)

Convert to a matrix in device memory.

**Parameters columns: sequence of str**

List of a column names to be extracted. The order is preserved. If `None` is specified, all columns are used.

**Returns** A (nrow x ncol) numpy ndarray in “F” order.

**as\_matrix** (*columns=None*)

Covert to a matrix in host memory.

**Parameters columns: sequence of str**

List of a column names to be extracted. The order is preserved. If `None` is specified, all columns are used.

**Returns** A (nrow x ncol) numpy ndarray in “F” order.

**columns**

Returns a tuple of columns

**copy** ()

Shallow copy this dataframe

**drop\_column** (*name*)

Drop a column by *name*

**dtypes**

Return the dtypes in this object.

**classmethod from\_pandas** (*dataframe*)

Convert from a Pandas DataFrame.

**Raises** `TypeError` for invalid input type.

**classmethod from\_records** (*data, index=None, columns=None*)

Convert from a numpy recarray or structured array.

**Parameters data** : numpy structured dtype or recarray

**index** : str

The name of the index column in *data*. If `None`, the default index is used.

**columns** : list of str

List of column names to include.

**Returns** DataFrame

**groupby** (*by*, *sort=False*, *as\_index=False*)  
Groupby

**Parameters** *by* : list-of-str or str

Column name(s) to form that groups by.

**sort** : bool

Force sorting group keys. Depends on the underlying algorithm. Current algorithm always sort.

**as\_index** : bool; defaults to False

Must be False. Provided to be API compatible with pandas. The keys are always left as regular columns in the result.

**Returns** The groupby object

## Notes

Unlike pandas, this groupby operation behaves like a SQL groupby. No empty rows are returned. (For categorical keys, pandas returns rows for all categories even if they are no corresponding values.)

Only a minimal number of operations is implemented so far.

- Only *by* argument is supported.
- The output is always sorted according to the *by* columns.
- Since we don't support multiindex, the *by* columns are stored as regular columns.

## index

Returns the index of the DataFrame

**join** (*other*, *on=None*, *how='left'*, *lsuffix=""*, *rsuffix=""*, *sort=False*)

Join columns with other DataFrame on index or on a key column.

**Parameters** *other* : DataFrame

**how** : str

Only accepts "left", "right", "inner", "outer"

**lsuffix, rsuffix** : str

The suffices to add to the left (*lsuffix*) and right (*rsuffix*) column names when avoiding conflicts.

**sort** : bool

Set to True to ensure sorted ordering.

**Returns** *joined* : DataFrame

## Notes

Difference from pandas:

- *other* must be a single DataFrame for now.
- *on* is not supported yet due to lack of multi-index support.

**label\_encoding** (*column*, *prefix*, *cats*, *prefix\_sep*='\_', *dtype*=None, *na\_sentinel*=-1)  
 Encode labels in a column with label encoding.

**Parameters** *column* : str

the source column with binary encoding for the data.

**prefix** : str

the new column name prefix.

**cats** : sequence of ints

the sequence of categories as integers.

**prefix\_sep** : str

the separator between the prefix and the category.

**dtype** :

the dtype for the outputs; see Series.label\_encoding

**na\_sentinel** : number

Value to indicate missing category.

**Returns**

—————  
 a new dataframe with a new column append for the coded values.

**loc**

Returns a label-based indexer for row-slicing and column selection.

## Examples

```
>>> df = DataFrame([('a', list(range(20))),
...                 ('b', list(range(20))),
...                 ('c', list(range(20)))]
# get rows from index 2 to index 5 from 'a' and 'b' columns.
>>> df.loc[2:5, ['a', 'b']]
   a  b
2  2  2
3  3  3
4  4  4
5  5  5
```

**nlargest** (*n*, *columns*, *keep*='first')

Get the rows of the DataFrame sorted by the *n* largest value of *columns*

Difference from pandas: \* Only a single column is supported in *columns*

**nsmallest** (*n*, *columns*, *keep*='first')

Get the rows of the DataFrame sorted by the *n* smallest value of *columns*

Difference from pandas: \* Only a single column is supported in *columns*

**one\_hot\_encoding** (*column*, *prefix*, *cats*, *prefix\_sep*='\_', *dtype*='float64')

Expand a column with one-hot-encoding.

**Parameters** *column* : str

the source column with binary encoding for the data.

**prefix** : str

the new column name prefix.

**cats** : sequence of ints

the sequence of categories as integers.

**prefix\_sep** : str

the separator between the prefix and the category.

**dtype** :

the dtype for the outputs; defaults to float64.

**Returns** a new dataframe with new columns append for each category.

**query** (*expr*)

Query with a boolean expression using Numba to compile a GPU kernel.

See pandas.DataFrame.query.

**Parameters** **expr** : str

A boolean expression. Names in the expression refers to the columns. Any name prefixed with @ refer to the variables in the calling environment.

**Returns** **filtered** : DataFrame

**set\_index** (*index*)

Return a new DataFrame with a new index

**Parameters** **index** : Index, Series-convertible or str

Index: the new index. Series-convertible: values for the new index. str: name of column to be used as series

**sort\_index** (*ascending=True*)

Sort by the index

**sort\_values** (*by, ascending=True*)

Sort by values.

Difference from pandas: \* *by* must be the name of a single column. \* Support axis='index' only. \* Not supporting: inplace, kind, na\_position

Details: Uses parallel radixsort, which is a stable sort.

**to\_pandas** ()

Convert to a Pandas DataFrame.

**to\_records** (*index=True*)

Covert to a numpy recarray

**Parameters** **index** : bool

Whether to include the index in the output.

**Returns** numpy recarray

**to\_string** (*nrows=NOTSET, ncols=NOTSET*)

Convert to string

**Parameters** **nrows** : int

Maximum number of rows to show. If it is None, all rows are shown.

**ncols** : int

Maximum number of columns to show. If it is None, all columns are shown.

## 1.2 Series

**class** pygdf.dataframe.**Series** (*data, index=None*)

Data and null-masks.

Series objects are used as columns of DataFrame.

### Attributes

<code>cat</code>	
<code>data</code>	The gpu buffer for the data
<code>dtype</code>	dtype of the Series
<code>has_null_mask</code>	A boolean indicating whether a null-mask is needed
<code>index</code>	The index object
<code>null_count</code>	Number of null values
<code>nullmask</code>	The gpu buffer for the null-mask
<code>valid_count</code>	Number of non-null values

### Methods

<code>append(arbitrary)</code>	Append values from another <code>Series</code> or array-like object.
<code>applymap(udf[, out_dtype])</code>	Apply a elementwise function to transform the values in the Column.
<code>argsort([ascending])</code>	Returns a Series of int64 index that will sort the series.
<code>as_index()</code>	
<code>as_mask()</code>	Convert booleans to bitmask
<code>astype(dtype)</code>	Convert to the given <code>dtype</code> .
<code>ceil()</code>	Rounds each value upward to the smallest integral value not less than the original.
<code>count()</code>	The number of non-null values
<code>factorize([na_sentinel])</code>	Encode the input values as integer labels
<code>fillna(value)</code>	Fill null values with <code>value</code> .
<code>find_first_value(value)</code>	Returns offset of first value that matches
<code>find_last_value(value)</code>	Returns offset of last value that matches
<code>floor()</code>	Rounds each value downward to the largest integral value not greater than the original.
<code>from_categorical(categorical[, codes])</code>	Creates from a <code>pandas.Categorical</code>
<code>from_masked_array(data, mask[, null_count])</code>	Create a Series with null-mask.
<code>head([n])</code>	
<code>label_encoding(cats[, dtype, na_sentinel])</code>	Perform label encoding
<code>max()</code>	Compute the max of the series
<code>mean()</code>	Compute the mean of the series

Continued on next page

Table 1.4 – continued from previous page

<code>mean_var([ddof])</code>	Compute mean and variance at the same time.
<code>min()</code>	Compute the min of the series
<code>nlargest([n, keep])</code>	Returns a new Series of the <i>n</i> largest element.
<code>nsmallest([n, keep])</code>	Returns a new Series of the <i>n</i> smallest element.
<code>one_hot_encoding(cats[, dtype])</code>	Perform one-hot-encoding
<code>reset_index()</code>	Reset index to RangeIndex
<code>reverse()</code>	Reverse the Series
<code>scale()</code>	Scale values to [0, 1] in float64
<code>set_index(index)</code>	Returns a new Series with a different index.
<code>set_mask(mask[, null_count])</code>	Create new Series by setting a mask array.
<code>sort_index([ascending])</code>	Sort by the index.
<code>sort_values([ascending])</code>	Sort by values.
<code>std([ddof])</code>	Compute the standard deviation of the series
<code>sum()</code>	Compute the sum of the series
<code>take(indices[, ignore_index])</code>	Return Series by taking values from the corresponding <i>indices</i> .
<code>to_array([fillna])</code>	Get a dense numpy array for the data.
<code>to_gpu_array([fillna])</code>	Get a dense numba device array for the data.
<code>to_pandas([index])</code>	
<code>to_string([nrows])</code>	Convert to string
<code>unique()</code>	Returns unique values of this Series.
<code>unique_k(k)</code>	
<code>values_to_string([nrows])</code>	Returns a list of string for each element.
<code>var([ddof])</code>	Compute the variance of the series

**append** (*arbitrary*)

Append values from another Series or array-like object. Returns a new copy with the index resetted.

**applymap** (*udf, out\_dtype=None*)

Apply a elementwise function to transform the values in the Column.

The user function is expected to take one argument and return the result, which will be stored to the output Series. The function cannot reference globals except for other simple scalar objects.

**Parameters** *udf* : function

Wrapped by `numba.cuda.jit` for call on the GPU as a device function.

**out\_dtype** : `numpy.dtype`; optional

The dtype for use in the output. By default, the result will have the same dtype as the source.

**Returns** *result* : Series

The mask and index are preserved.

**argsort** (*ascending=True*)

Returns a Series of int64 index that will sort the series.

Uses stable parallel radixsort.

**Returns** *result*: Series

**as\_mask** ()

Convert booleans to bitmask

**Returns** device array

**astype** (*dtype*)

Convert to the given dtype.

**Returns** If the dtype changed, a new `Series` is returned by casting each values to the given dtype.

If the dtype is not changed, `self` is returned.

**ceil** ()

Rounds each value upward to the smallest integral value not less than the original.

Returns a new `Series`.

**count** ()

The number of non-null values

**data**

The gpu buffer for the data

**dtype**

dtype of the `Series`

**factorize** (*na\_sentinel=-1*)

Encode the input values as integer labels

**Parameters** `na_sentinel` : number

Value to indicate missing category.

**Returns** (`labels`, `cats`) : (`Series`, `Series`)

- `labels` contains the encoded values
- `cats` contains the categories in order that the N-th item corresponds to the (N-1) code.

**fillna** (*value*)

Fill null values with `value`.

Returns a copy with null filled.

**find\_first\_value** (*value*)

Returns offset of first value that matches

**find\_last\_value** (*value*)

Returns offset of last value that matches

**floor** ()

Rounds each value downward to the largest integral value not greater than the original.

Returns a new `Series`.

**classmethod from\_categorical** (*categorical*, *codes=None*)

Creates from a `pandas.Categorical`

If `codes` is defined, use it instead of `categorical.codes`

**classmethod from\_masked\_array** (*data*, *mask*, *null\_count=None*)

Create a `Series` with null-mask. This is equivalent to:

```
Series(data).set_mask(mask, null_count=null_count)
```

**Parameters** `data` : 1D array-like

The values. Null values must not be skipped. They can appear as garbage values.

**mask** : 1D array-like of `numpy.uint8`

The null-mask. Valid values are marked as 1; otherwise 0. The mask bit given the data index `idx` is computed as:

```
(mask[idx // 8] >> (idx % 8)) & 1
```

**null\_count** : int, optional

The number of null values. If None, it is calculated automatically.

**has\_null\_mask**

A boolean indicating whether a null-mask is needed

**index**

The index object

**label\_encoding** (*cats, dtype=None, na\_sentinel=-1*)

Perform label encoding

**Parameters values** : sequence of input values

**dtype: numpy.dtype; optional**

Specifies the output dtype. If *None* is given, the smallest possible integer dtype (starting with `np.int32`) is used.

**na\_sentinel** : number

Value to indicate missing category.

**Returns**

\_\_\_\_\_

**A sequence of encoded labels with value between 0 and n-1 classes(cats)**

**max** ()

Compute the max of the series

**mean** ()

Compute the mean of the series

**mean\_var** (*ddof=1*)

Compute mean and variance at the same time.

**min** ()

Compute the min of the series

**nlargest** (*n=5, keep='first'*)

Returns a new Series of the *n* largest element.

**nsmallest** (*n=5, keep='first'*)

Returns a new Series of the *n* smallest element.

**null\_count**

Number of null values

**nullmask**

The gpu buffer for the null-mask

**one\_hot\_encoding** (*cats, dtype='float64'*)

Perform one-hot-encoding

**Parameters cats** : sequence of values

values representing each category.



**dtype** : numpy.dtype  
 specifies the output dtype.

**Returns** A sequence of new series for each category. Its length is determined by the length of `cats`.

**reset\_index()**  
 Reset index to RangeIndex

**reverse()**  
 Reverse the Series

**scale()**  
 Scale values to [0, 1] in float64

**set\_index(index)**  
 Returns a new Series with a different index.

**Parameters index** : Index, Series-convertible  
 the new index or values for the new index

**set\_mask(mask, null\_count=None)**  
 Create new Series by setting a mask array.

This will override the existing mask. The returned Series will reference the same data buffer as this Series.

**Parameters mask** : 1D array-like of numpy.uint8

The null-mask. Valid values are marked as 1; otherwise 0. The mask bit given the data index `idx` is computed as:

```
(mask[idx // 8] >> (idx % 8)) & 1
```

**null\_count** : int, optional

The number of null values. If None, it is calculated automatically.

**sort\_index(ascending=True)**  
 Sort by the index.

**sort\_values(ascending=True)**  
 Sort by values.

Difference from pandas: \* Support axis='index' only. \* Not supporting: inplace, kind, na\_position

Details: Uses parallel radixsort, which is a stable sort.

**std(ddof=1)**  
 Compute the standard deviation of the series

**sum()**  
 Compute the sum of the series

**take(indices, ignore\_index=False)**  
 Return Series by taking values from the corresponding *indices*.

**to\_array(fillna=None)**  
 Get a dense numpy array for the data.

**Parameters fillna** : str or None

Defaults to None, which will skip null values. If it equals "pandas", null values are filled with NaNs. Non integral dtype is promoted to np.float64.

### Notes

if `fillna` is `None`, null values are skipped. Therefore, the output size could be smaller.

**to\_gpu\_array** (*fillna=None*)

Get a dense numba device array for the data.

**Parameters** `fillna` : str or None

See `fillna` in `.to_array`.

### Notes

if `fillna` is `None`, null values are skipped. Therefore, the output size could be smaller.

**to\_string** (*nrows=NOTSET*)

Convert to string

**Parameters** `nrows` : int

Maximum number of rows to show. If it is `None`, all rows are shown.

**unique** ()

Returns unique values of this Series.

**valid\_count**

Number of non-null values

**values\_to\_string** (*nrows=None*)

Returns a list of string for each element.

**var** (*ddof=1*)

Compute the variance of the series

## 1.3 Groupby

**class** `pygdf.groupby.Groupby` (*df, by*)

Groupby object returned by `pygdf.DataFrame.groupby()`.

### Methods

<code>agg(args)</code>	Invoke aggregation functions on the groups.
<code>apply(function)</code>	Apply a transformation function over the grouped chunk.
<code>apply_grouped(function, **kwargs)</code>	
<code>as_df()</code>	Get the intermediate dataframe after shuffling the rows into groups.
<code>count()</code>	Compute the count of each group
<code>max()</code>	Compute the max of each group
<code>mean()</code>	Compute the mean of each group
<code>min()</code>	Compute the min of each group
<code>std()</code>	Compute the std of each group

**agg** (*args*)

Invoke aggregation functions on the groups.

**Parameters** **args**: **dict**, **list**, **str**, **callable**

- **str** The aggregate function name.
- **callable** The aggregate function.
- **list** List of *str* or *callable* of the aggregate function.
- **dict** key-value pairs of source column name and list of aggregate functions as *str* or *callable*.

**Returns** **result** : DataFrame**apply** (*function*)

Apply a transformation function over the grouped chunk.

**as\_df** ()

Get the intermediate dataframe after shuffling the rows into groups.

**Returns** (**df**, **segs**) : namedtuple

- **df** : DataFrame
- **segs** [Series] Beginning offsets of each group.

**count** ()

Compute the count of each group

**Returns** **result** : DataFrame**max** ()

Compute the max of each group

**Returns** **result** : DataFrame**mean** ()

Compute the mean of each group

**Returns** **result** : DataFrame**min** ()

Compute the min of each group

**Returns** **result** : DataFrame**std** ()

Compute the std of each group

**Returns** **result** : DataFrame

## 1.4 GpuArrowReader

**class** pygdf.gpuarrow.**GpuArrowReader** (*schema\_data*, *gpu\_data*)

### Methods

---

*to\_dict*()Return a dictionary of Series object

---

`to_dict()`

Return a dictionary of Series object

## CHAPTER 2

---

### Indices and tables

---

- `genindex`
- `modindex`
- `search`



**A**

add\_column() (pygdf.dataframe.DataFrame method), 3  
 agg() (pygdf.groupby.Groupby method), 15  
 append() (pygdf.dataframe.Series method), 10  
 apply() (pygdf.groupby.Groupby method), 15  
 apply\_chunks() (pygdf.dataframe.DataFrame method), 3  
 apply\_rows() (pygdf.dataframe.DataFrame method), 4  
 applymap() (pygdf.dataframe.Series method), 10  
 argsort() (pygdf.dataframe.Series method), 10  
 as\_df() (pygdf.groupby.Groupby method), 15  
 as\_gpu\_matrix() (pygdf.dataframe.DataFrame method), 5  
 as\_mask() (pygdf.dataframe.Series method), 10  
 as\_matrix() (pygdf.dataframe.DataFrame method), 5  
 astype() (pygdf.dataframe.Series method), 10

**C**

ceil() (pygdf.dataframe.Series method), 11  
 columns (pygdf.dataframe.DataFrame attribute), 5  
 copy() (pygdf.dataframe.DataFrame method), 5  
 count() (pygdf.dataframe.Series method), 11  
 count() (pygdf.groupby.Groupby method), 15

**D**

data (pygdf.dataframe.Series attribute), 11  
 DataFrame (class in pygdf.dataframe), 1  
 drop\_column() (pygdf.dataframe.DataFrame method), 5  
 dtype (pygdf.dataframe.Series attribute), 11  
 dtypes (pygdf.dataframe.DataFrame attribute), 5

**F**

factorize() (pygdf.dataframe.Series method), 11  
 fillna() (pygdf.dataframe.Series method), 11  
 find\_first\_value() (pygdf.dataframe.Series method), 11  
 find\_last\_value() (pygdf.dataframe.Series method), 11  
 floor() (pygdf.dataframe.Series method), 11  
 from\_categorical() (pygdf.dataframe.Series class method), 11  
 from\_masked\_array() (pygdf.dataframe.Series class method), 11

from\_pandas() (pygdf.dataframe.DataFrame class method), 5  
 from\_records() (pygdf.dataframe.DataFrame class method), 5

**G**

GpuArrowReader (class in pygdf.gpuarrow), 15  
 Groupby (class in pygdf.groupby), 14  
 groupby() (pygdf.dataframe.DataFrame method), 5

**H**

has\_null\_mask (pygdf.dataframe.Series attribute), 12

**I**

index (pygdf.dataframe.DataFrame attribute), 6  
 index (pygdf.dataframe.Series attribute), 12

**J**

join() (pygdf.dataframe.DataFrame method), 6

**L**

label\_encoding() (pygdf.dataframe.DataFrame method), 6  
 label\_encoding() (pygdf.dataframe.Series method), 12  
 loc (pygdf.dataframe.DataFrame attribute), 7

**M**

max() (pygdf.dataframe.Series method), 12  
 max() (pygdf.groupby.Groupby method), 15  
 mean() (pygdf.dataframe.Series method), 12  
 mean() (pygdf.groupby.Groupby method), 15  
 mean\_var() (pygdf.dataframe.Series method), 12  
 min() (pygdf.dataframe.Series method), 12  
 min() (pygdf.groupby.Groupby method), 15

**N**

nlargest() (pygdf.dataframe.DataFrame method), 7  
 nlargest() (pygdf.dataframe.Series method), 12  
 nsmallest() (pygdf.dataframe.DataFrame method), 7

nsmallest() (pygdf.dataframe.Series method), 12  
null\_count (pygdf.dataframe.Series attribute), 12  
nullmask (pygdf.dataframe.Series attribute), 12

## O

one\_hot\_encoding() (pygdf.dataframe.DataFrame method), 7  
one\_hot\_encoding() (pygdf.dataframe.Series method), 12

## Q

query() (pygdf.dataframe.DataFrame method), 8

## R

reset\_index() (pygdf.dataframe.Series method), 13  
reverse() (pygdf.dataframe.Series method), 13

## S

scale() (pygdf.dataframe.Series method), 13  
Series (class in pygdf.dataframe), 9  
set\_index() (pygdf.dataframe.DataFrame method), 8  
set\_index() (pygdf.dataframe.Series method), 13  
set\_mask() (pygdf.dataframe.Series method), 13  
sort\_index() (pygdf.dataframe.DataFrame method), 8  
sort\_index() (pygdf.dataframe.Series method), 13  
sort\_values() (pygdf.dataframe.DataFrame method), 8  
sort\_values() (pygdf.dataframe.Series method), 13  
std() (pygdf.dataframe.Series method), 13  
std() (pygdf.groupby.Groupby method), 15  
sum() (pygdf.dataframe.Series method), 13

## T

take() (pygdf.dataframe.Series method), 13  
to\_array() (pygdf.dataframe.Series method), 13  
to\_dict() (pygdf.gpuarrow.GpuArrowReader method), 16  
to\_gpu\_array() (pygdf.dataframe.Series method), 14  
to\_pandas() (pygdf.dataframe.DataFrame method), 8  
to\_records() (pygdf.dataframe.DataFrame method), 8  
to\_string() (pygdf.dataframe.DataFrame method), 8  
to\_string() (pygdf.dataframe.Series method), 14

## U

unique() (pygdf.dataframe.Series method), 14

## V

valid\_count (pygdf.dataframe.Series attribute), 14  
values\_to\_string() (pygdf.dataframe.Series method), 14  
var() (pygdf.dataframe.Series method), 14