

---

# **Kodexa**

***Release 2.0.17***

**May 29, 2020**



---

## Contents:

---

<b>1</b>	<b>Kodexa</b>	<b>3</b>
<b>2</b>	<b>Mix-Ins</b>	<b>5</b>
2.1	Core . . . . .	5
2.2	Spatial . . . . .	10
<b>3</b>	<b>API Documentation</b>	<b>13</b>
3.1	Kodexa Cloud . . . . .	13
3.2	Connectors . . . . .	13
3.3	Core Model . . . . .	13
3.4	Pipeline . . . . .	17
3.5	Sinks . . . . .	19
3.6	Stores . . . . .	20
<b>4</b>	<b>Indices and tables</b>	<b>23</b>
	<b>Python Module Index</b>	<b>25</b>
	<b>Index</b>	<b>27</b>



This the documentation for 2.0.17



# CHAPTER 1

---

## Kodexa

---

Kodexa is a core python package for the Kodexa Cloud, and a surrounding set of extensions, that are built to allow developers to work with documents and unstructured data in a way that supports flexibility, extensibility and deployable in a range of platforms.





In this section we cover some of the included mix-ins and detail the functions that they add.

## 2.1 Core

It is the core set of helpers that can be added to content nodes to allow for searching and tagging.

**tag**(*self*, *tag\_name*, *type\_re=None*, *content\_re=None*, *use\_all\_content=False*, *node\_only=False*, *include\_children=False*, *fixed\_position=None*, *data=None*)

This will tag (see Feature Tagging) the expression groups identified by the regular expression.

```
>>> document.content_node.find(content_re='.*Cheese.*').tag('is_cheese')
```

### Parameters

- **tag\_name** – the name of the tag to apply
- **type\_re** – regular expression to make the type (default `.*`)
- **content\_re** – the regular expression that you wish to use to tag, note that we will create a tag for each matching group
- **use\_all\_content** – apply the regular expression to the all\_content (include content from child nodes)
- **node\_only** – Ignore the matching groups and tag the whole node
- **include\_children** – Include recurse into children and tag where matching
- **fixed\_position** – use a fixed position, supplied as a tuple i.e. - (4,10) tag from position 4 to 10 (default None)
- **data** – Attach the a dictionary of data for the given tag

**tag\_range**(*self*, *start\_content\_re*, *end\_content\_re*, *tag\_name*, *type\_re='.\*'*, *use\_all\_content=False*)

This will tag all the child nodes between the start and end content regular expressions

```
>>> document.content_node.tag_range(start_content_re='.*Cheese.*', end_content_re=
↳'.*Fish.*', tag_name='foo')
```

### Parameters

- **start\_content\_re** – The regular expression to match the starting child
- **end\_content\_re** – The regular expression to match the ending child
- **tag\_name** – The tag to be applied to the nodes in range
- **type\_re** – The type to match (default is all)
- **use\_all\_content** – Use full content (including child nodes, default is False)

**get\_all\_content** (*self*, *separator*='')

This will build the complete content, including the content of children.

```
>>> document.content_node.get_all_content()
"This string is made up of multiple nodes"
```

**Parameters separator** – the separate to use to join the content together (default is " ")

**get\_all\_tags** (*self*)

Returns a list of the names of the tags on the given node and all its children

```
>>> document.content_node.find(content_re='.*Cheese.*').get_all_tags()
['is_cheese']
```

**Returns** A list of the tag names

**get\_tag** (*self*, *tag\_name*)

Returns the value of a tag, this can be either a single list [start,end,value] or if multiple parts of the content of this node match you can end up with a list of lists i.e. [[start1,end1,value1],[start2,end2,value2]]

```
>>> document.content_node.find(content_re='.*Cheese.*').get_tag('is_cheese')
[0,10, 'The Cheese Moved']
```

**Parameters tag\_name** – The name of the tag

**Returns** The tagged location and value (or a list if more than one)

**find** (*self*, *content\_re*='.\*', *type\_re*='.\*', *direction*=<FindDirection.CHILDREN: 1>, *tag\_name*=None, *instance*=0, *tag\_name\_re*=None, *use\_all\_content*=False)

Search for a node that matches on the value and or type using regular expressions

```
>>> document.content_node.find(content_re='.*Cheese.*', instance=2)
<kodexa.model.model.ContentNode object at 0x7f80605e53c8>
```

### Parameters

- **content\_re** – the regular expression to match the nodes content (default '.\*')
- **type\_re** – the regular expression to match the nodes type (default '.\*')
- **direction** – the direction to search, either FindDirection.CHILDREN or FindDirection.PARENT (default CHILDREN)

- **tag\_name** – the tag name that must exist
- **instance** – the instance to return (0=first)
- **tag\_name\_re** – the regular expression to match for the tag\_name
- **use\_all\_content** – match the content for all child nodes to (default False)

**Returns** either an instance of ContentNode (if found), or None

```
findall (self, content_re='.*', type_re='.*', direction=<FindDirection.CHILDREN: 1>, tag_name=None,
         tag_name_re=None, use_all_content=False)
```

Search for nodes that matches on the value and or type using regular expressions

```
>>> document.content_node.findall(content_re='.*Cheese.*')
[<kodexa.model.model.ContentNode object at 0x7f80605e53c8>,
 <kodexa.model.model.ContentNode object at 0x7f80605e53c8>]
```

### Parameters

- **content\_re** – the regular expression to match the nodes content (default '.\*')
- **type\_re** – the regular expression to match the nodes type (default '.\*')
- **direction** – the direction to search, either FindDirection.CHILDREN or FindDirection.PARENT (default CHILDREN)
- **tag\_name** – the tag name that must exist
- **tag\_name\_re** – the regular expression to match for the tag\_name
- **use\_all\_content** – match the content for all child nodes to (default False)

**Returns** list of matching content nodes

```
find_with_feature_value (self, feature_type, feature, value, direction=<FindDirection.CHILDREN:
                          1>, instance=1)
```

Search for a node with a specific feature type, name and value

```
>>> document.content_node.find_with_feature_value(feature_type='tag', feature='is_
↪cheese', value=[1,10,'The Cheese has moved'])
<kodexa.model.model.ContentNode object at 0x7f80605e53c8>
```

### Parameters

- **feature\_type** – the feature type
- **feature** – the feature name
- **value** – the feature value
- **direction** – the direction to search, either FindDirection.CHILDREN or FindDirection.PARENT (default CHILDREN)
- **instance** – the instance to get (default 1)

**Returns** either an instance of ContentNode (if found), or None

```
findall_with_feature_value (self, feature_type, feature_name, value, direc-
                             tion=<FindDirection.CHILDREN: 1>)
```

Search for all nodes with a specific feature type, name and value

```
>>> document.content_node.findall_with_feature_value(feature_type='tag', feature=
↳ 'is_cheese', value=[1,10, 'The Cheese has moved'])
[<kodexa.model.model.ContentNode object at 0x7f80605e53c8>]
```

#### Parameters

- **feature\_type** – the feature type
- **feature\_name** – the feature name
- **value** – the feature value
- **direction** – the direction to search, either FindDirection.CHILDREN or FindDirection.PARENT (default CHILDREN)

**Returns** list of the matching content nodes

#### **get\_tags** (*self*)

Returns a list of the names of the tags on the given node

```
>>> document.content_node.find(content_re='.*Cheese.*').get_tags()
['is_cheese']
```

**Returns** A list of the tag name

#### **move\_child\_to\_parent** (*self, target\_child, target\_parent*)

This will move the *target\_child*, which must be a child of the node, to a new parent.

It will be added to the end of the parent

```
>>> document.content_root.move_child_to_parent(document.content_root.find(type_
↳ ref='line'), document.content_root)
```

#### Parameters

- **target\_child** – the child node that needs to be moved
- **target\_parent** – the parent to attach this node to

#### **adopt\_children** (*self, children, replace=False*)

This will take a list of content nodes and adopt them under this node, ensuring they are re-parented.

It will be added to the end of the parent

```
>>> document.content_root.adopt_children(document.content_root.find(type_ref='line
↳ '), replace=True)
```

#### Parameters

- **children** – a list of the children to adopt
- **replace** – if True the node will remove all current children and replace them with the new list

#### **remove\_tag** (*self, tag\_name*)

This will remove a tag from the given node

```
>>> document.content_node.remove_tag(tag_name='foo')
```

**Parameters** `tag_name` – The tag to be applied to the nodes in range

**collect\_nodes\_to** (*self*, *end\_node*)

Return a list of the sibling nodes between the current node and the end node

```
>>> document.content_node.children[0].collect_nodes_to(end_node=document.content_
↳node.children[5])
```

**Parameters** `end_node` – The node to end at

**tag\_nodes\_to** (*self*, *end\_node*, *tag\_name*)

Tag all the nodes from this node to the end node with the given tag name

```
>>> document.content_node.children[0].tag_nodes_to(document.content_node.
↳children[5], tag_name='foo')
```

**param** `end_node` The node to end with

**param** `tag_name` The tag name

**get\_node\_at\_index** (*self*, *index*)

Returns the node at a specific index, if the index is outside the first (0), or last index it will return null, if not it will return the node at index, if there isn't a node at that index it will return a 'virtual' node that will represent the node to the side of this node and have an index and no features or content

**Parameters** `index` – The index (zero-based) for the child node

**Returns** Node at index, or None if the index is outside the boundaries of child nodes

**has\_next\_node** (*self*, *type\_re*='.\*', *skip\_virtual*=False)

Returns True if the node has a next node

**Parameters**

- **type\_re** – Type name (regular expression)
- **skip\_virtual** – True to skip any virtual nodes and only return the next real node (default False)

**Returns** True if there is a next sibling node

**has\_previous\_node** (*self*, *type\_re*='.\*', *skip\_virtual*=False)

Returns True if the node has a previous node

**Parameters**

- **type\_re** – Type name (regular expression)
- **skip\_virtual** – True to skip any virtual nodes and only return the next real node (default False)

**return** True if there is a previous sibling node

**next\_node** (*self*, *type\_re*='.\*', *skip\_virtual*=False, *has\_no\_content*=False, *traverse*=<Traverse.SIBLING: 1>)

Returns the next sibling content node, note that this logic is based on the index, therefore the next node might actually be a virtual node that is created to fill a gap in the document, since the index allows for sparse documents

#### Parameters

- **type\_re** – the regular expression for the type of node (default ‘.\*’)
- **skip\_virtual** – True to skip any virtual nodes and only return the next real node (default False)
- **has\_no\_content** – True if you only want to return a node that has no content
- **traverse** – By default we traverse siblings, however you can include CHILDREN, PARENT or ALL

**Returns** the next node or None if no node exists

**previous\_node** (*self*, *type\_re*='.\*', *skip\_virtual*=False, *has\_no\_content*=False, *traverse*=<Traverse.SIBLING: 1>)

Returns the previous sibling content node, note that this logic is based on the index, therefore the next node might actually be a virtual node that is created to fill a gap in the document, since the index allows for sparse documents

#### Parameters

- **type\_re** – the regular expression for the type of node (default ‘.\*’)
- **skip\_virtual** – True to skip any virtual nodes and only return the next real node (default False)
- **has\_no\_content** – True if you only want to return a node that no content
- **traverse** – By default we traverse siblings, however you can include CHILDREN, PARENT or ALL

**Returns** the previous node or None if no node exists

**get\_last\_child\_index** (*self*)

Returns the max index value for the children of this node, if the node has no children it returns None

**Returns** The max index of the children of this node, or None if no children

**is\_first\_child** (*self*)

Returns True if this is the first child, also True if it has no parent

**Returns** True if this is the first child

**is\_last\_child** (*self*)

Returns True if this is the last child, also True if it has no parent

**Returns** True if this is the first child

## 2.2 Spatial

One of the core mix-ins is Spatial. It is based on the concept of holding spatial information about the content nodes.

This spatial information can then be used by the mix-in’s methods to allow you to both pull spatial information, but also to query it.

**set\_statistics** (*self*, *statistics*)

Set the spatial statistics for this node

```
>>> document.content_node.find(type_re='page').set_statistics(NodeStatistics())
```

**Parameters** `statistics` – the statistics object

**get\_statistics** (*self*)

Get the spatial statistics for this node

```
>>> document.content_node.find(type_re='page').get_statistics()
<kodexa.spatial.NodeStatistics object at 0x7f80605e53c8>
```

**Returns** the statistics object (or None if not set)

**set\_bbox** (*self*, *bbox*)

Set the bounding box for the node, this is structured as:

[x1,y1,x2,y2]

```
>>> document.content_node.find(type_re='page').set_bbox([10,20,50,100])
```

**Parameters** `bbox` – the bounding box array

**get\_bbox** (*self*)

Get the bounding box for the node, this is structured as:

[x1,y1,x2,y2]

```
>>> document.content_node.find(type_re='page').get_bbox()
[10,20,50,100]
```

**Returns** the bounding box array

**set\_rotate** (*self*, *rotate*)

Set the rotate of the node

```
>>> document.content_node.find(type_re='page').set_rotate(90)
```

:param rotate the rotation of the node

**get\_rotate** (*self*)

Get the rotate of the node

```
>>> document.content_node.find(type_re='page').get_rotate()
90
```

**Returns** the rotation of the node

**get\_x** (*self*)

Get the X position of the node

```
>>> document.content_node.find(type_re='page').get_x()
10
```

**Returns** the X position of the node

**get\_y** (*self*)

Get the Y position of the node

```
>>> document.content_node.find(type_re='page').get_y()  
90
```

**Returns** the Y position of the node

**get\_width** (*self*)

Get the width of the node

```
>>> document.content_node.find(type_re='page').get_width()  
70
```

**Returns** the width of the node

**get\_height** (*self*)

Get the height of the node

```
>>> document.content_node.find(type_re='page').get_height()  
40
```

**Returns** the height of the node

**set\_bbox\_from\_children** (*self*)

Set the bounding box for this node based on its children

**collapse** (*self, type\_re*)

Will collapse the given type, this will remove this type from the hierarchy.

**Parameters** **type\_re** – the type that you will collapse

**Returns**



Kodexa is a Python framework to enable flexible data engineering with semi-structured and unstructured documents and data.

## 3.1 Kodexa Cloud

Provides out of the box integration with the Kodexa cloud, enabling the universe of content services that are available

```
class KodexaCloudPipeline (slug, version=None, attach_source=True, options=None, auth=None,  
                           cloud_url='https://cloud.kodexa.com', access_token=None)
```

Allow you to interact with a pipeline that has been deployed in the Kodexa Cloud

```
class KodexaCloudService (slug, version=None, attach_source=False, options={}, auth=[],  
                           cloud_url='https://cloud.kodexa.com', access_token=None)
```

Allows you to interact with a content service that has been deployed in the Kodexa cloud

## 3.2 Connectors

Connectors provide a way to access document (files or otherwise) from a source, and they form the starting point for Pipelines

```
class FolderConnector (path, file_filter='*')
```

```
class UrlConnector (url, headers=None)
```

## 3.3 Core Model

The core model provides the object structure for Documents, ContentNodes and Features which is used as the foundation for working with unstructured data in the framework.

Create a new instance of a Document, you will be required to provide a DocumentMetadata object

```
>>> document = Document(DocumentMetadata())
```

```
class Document (metadata=None, content_node: kodexa.model.model.ContentNode = None,
                 source=<kodexa.model.model.SourceMetadata object>)
```

A Document is a collection of metadata and a set of content nodes.

```
add_mixin (mixin)
```

Add the given mixin to this document, this will apply the mixin to all the content nodes, and also register it with the document so that future invocations of `create_node` will ensure the node has the mixin applied.

```
>>> document.add_mixin('spatial')
```

```
create_node (type: str = <class 'type'>, content: str = None, virtual: bool = False, parent:
              kodexa.model.model.ContentNode = None, index: int = 0)
```

Creates a new node for the document, this doesn't add the node to the document however it does ensure that any mixins that have been applied to the document will also be available on the new node

```
>>> document.create_node(type='page')
<kodexa.model.model.ContentNode object at 0x7f80605e53c8>
```

```
static from_dict (doc_dict)
```

Build a new document from a dictionary

```
>>> Document.from_dict(doc_dict)
```

```
static from_json (json_string)
```

From a JSON string create an instance of a Document

```
>>> document.from_json()
```

```
static from_kdxa (file_path)
```

Read an .kdxa file from the given file\_path and

```
>>> document = Document.from_kdxa('my-document.kdxa')
```

**Parameters file\_path** – the path to the mdoc file

```
static from_msgpack (bytes)
```

From a message pack byte array create an instance of a Document

```
>>> document.from_msgpack()
```

```
get_mixins ()
```

Get the list of mixins that have been enabled on this document.

```
>>> document.get_mixins()
['spatial', 'finders']
```

```
get_root ()
```

Get the root content node for the document (same as `content_node`)

```
>>> node = document.get_node()
```

```
to_dict ()
```

Convert this document object structure into a simple set of dictionaries

```
>>> document.to_dict()
```

**to\_json()**

Convert this document object structure into a JSON object

```
>>> document.to_json()
```

**to\_kdxa** (*file\_path*)

Write the document to the kdxa format (msgpack) which can be used with the Kodexa platform

```
>>> document.to_mdoc('my-document.kdxa')
```

**Parameters** *file\_path* – the path to the mdoc you wish to create

**to\_msgpack()**

Convert this document object structure into a message pack

```
>>> document.to_msgpack()
```

**to\_text()**

Convert this document object structure into a text representation, which can be useful when trying to review the structure.

```
>>> document.to_text()
```

**class DocumentMetadata** (*\*args, \*\*kwargs*)

A flexible dict based approach to capturing metadata for the document

**class ContentNode** (*document, type, content=""*, *content\_parts=[]*)

A Content Node identifies a section of the document containing logical grouping of information

The node will have content and can include n number of features.

You should always create a node using the Document create\_node method to ensure that the correct mixins are applied.

```
>>> new_page = document.create_node(type='page')
<kodexa.model.model.ContentNode object at 0x7f80605e53c8>
>>> current_content_node.add_child(new_page)
```

or

```
>>> new_page = document.create_node(type='page', content='This is page 1')
<kodexa.model.model.ContentNode object at 0x7f80605e53c8>
>>> current_content_node.add_child(new_page)
```

**add\_child** (*child, index=None*)

Add a ContentNode as a child of this ContentNode

```
>>> new_page = document.create_node(type='page')
<kodexa.model.model.ContentNode object at 0x7f80605e53c8>
>>> current_content_node.add_child(new_page)
```

**add\_feature** (*feature\_type, name, value, single=True, serialized=False*)

Add a new feature to this ContentNode.

You will need to provide the feature type, the name of the feature and then the value.

Note this will add a value to an existing feature, therefore the feature value might switch to being a list

```
>>> new_page = document.create_node(type='page')
<kodexa.model.model.ContentNode object at 0x7f80605e53c8>
>>> new_page.add_feature('pagination', 'pageNum', 1)
```

**get\_children()**

Returns a list of the children of this node

```
>>> new_page = document.create_node(type='page')
<kodexa.model.model.ContentNode object at 0x7f80605e53c8>
>>> new_page.get_children()
[]
```

**get\_content()**

Returns the content of the node

```
>>> new_page.get_content()
"This is page one"
```

**get\_feature(feature\_type, name)**

Gets the value for the given feature.

You will need to provide the type and name of the feature. If no feature is found you will get None

```
>>> new_page.get_feature('pagination', 'pageNum')
1
```

**get\_feature\_value(feature\_type, name)**

Returns the assigned value for a given feature

You will need to provide the type and name of the feature. If the feature is present it will return the value otherwise it will return None.

```
>>> new_page.get_feature_value('pagination', 'pageNum')
1
```

**get\_features()**

Returns a list of the features on this content node

**Returns** a list of the features present

**get\_features\_of\_type(feature\_type)**

Return a list of all the features of a specific type

```
>>> new_page.get_features_of_type('tag')
[]
```

**Returns** list of the tags

**get\_type()**

Returns the type of the node

```
>>> new_page.get_content()
"page"
```

**has\_feature(feature\_type, name)**

Determines if the feature with the given name and type exists on this content node

You will need to provide the type and name of the feature. If the feature is present it will return True, else it will return False

```
>>> new_page.has_feature('pagination', 'pageNum')
True
```

**Returns** True if the feature is present

**remove\_feature** (*feature\_type, name*)

Determines if the feature with the given name and type exists on this content node

You will need to provide the type and name of the feature. If the feature is present it will return True, else it will return False

```
>>> new_page.remove_feature('pagination', 'pageNum')
```

**set\_feature** (*feature\_type, name, value*)

Sets a feature to this ContentNode, replacing the value

You will need to provide the feature type, the name of the feature and then the value.

Note this will replace any matching feature (i.e. with the same type and name)

```
>>> new_page = document.create_node(type='page')
<kodexa.model.model.ContentNode object at 0x7f80605e53c8>
>>> new_page.add_feature('pagination', 'pageNum', 1)
```

**to\_dict** ()

Convert the ContentNode, and all its children into a simple dictionary

```
>>> new_page = document.create_node(type='page')
<kodexa.model.model.ContentNode object at 0x7f80605e53c8>
>>> current_content_node.to_dict()
```

**to\_json** ()

Convert this node structure into a JSON object

```
>>> node.to_json()
```

**to\_text** ()

Convert this node structure into a text representation, which can be useful when trying to review the structure.

```
>>> node.to_text()
```

## 3.4 Pipeline

A Pipeline is a way to bring together a Connector, set of steps and then a sink to perform data cleansing, normalization, analysis and more.

**class Pipeline** (*connector, name='Default', stop\_on\_exception=True, logging\_level=20*)

A pipeline represents a way to bring together parts of the kodexa framework to solve a specific problem.

When you create a Pipeline you must provide the connector that will be used to source the documents.

```
>>> pipeline = Pipeline(FolderConnector(path='/tmp/', file_filter='example.pdf'))
```

### Parameters

- **connector** – the connector that will be the starting point for the pipeline
- **name** – the name of the pipeline (default 'Default')
- **stop\_on\_exception** – Should the pipeline raise exceptions and stop (default True)
- **logging\_level** – The logging level of the pipeline (default INFO)

### add\_step(*step*)

Add the given step to the current pipeline

```
>>> pipeline = Pipeline(FolderConnector(path='/tmp/', file_filter='example.pdf
↳'))
>>> pipeline.add_step(ExampleStep())
```

Note that it is also possible to add a function as a step, for example

```
>>> def my_function(doc):
>>>     doc.metadata.fishstick = 'foo'
>>>     return doc
>>> pipeline.add_step(my_function)
```

**Parameters** **step** – the step to add

### add\_store(*name, store*)

Add the store to the pipeline so that it is available to the pipeline

```
>>> pipeline = Pipeline(FolderConnector(path='/tmp/', file_filter='example.pdf
↳'))
>>> pipeline.add_store("test-store", InMemoryObjectStore())
```

### Parameters

- **name** – the name of the store (to refer to it)
- **store** – the store that should be added

### run()

Run the current pipeline, note that you must have a sink in place to allow the pipeline to run

```
>>> pipeline = Pipeline(FolderConnector(path='/tmp/', file_filter='example.pdf
↳'))
>>> pipeline.set_sink(ExampleSink())
>>> pipeline.run()
```

**Returns** The context from the run

### set\_sink(*sink*)

Set the sink you wish to use, note that it will replace any currently assigned sink

```
>>> pipeline = Pipeline(FolderConnector(path='/tmp/', file_filter='example.pdf
↳'))
>>> pipeline.set_sink(ExampleSink())
```

**Parameters** **sink** – the sink for the pipeline

```
class PipelineContext (content_provider=<kodexa.pipeline.pipeline.InMemoryContentProvider
                        object>, store_provider=<kodexa.pipeline.pipeline.InMemoryStoreProvider
                        object>, existing_content_objects=None, context=None)
```

Pipeline context is created when you create a pipeline and it provides a way to access information about the pipeline that is running. It can be made available to steps/functions so they can interact with it.

It also provides access to the ‘stores’ that have been added to the pipeline

```
add_store (name, store)
```

Add a store with given name to the context

**Parameters**

- **name** – the name to refer to the store with
- **store** – the instance of the store

```
get_store (name, default=None)
```

Get a store with given name from the context

**Parameters**

- **name** – the name to refer to the store with
- **default** – optionally the default to create the store as if it isn’t there

**Returns** the store, or None is not available

```
get_store_names ()
```

Return the list of store names in context

**Returns** the list of store names

```
set_output_document (output_document)
```

Set the output document from the pipeline

**Parameters** **output\_document** – the final output document from the pipeline

**Returns** the final output document

```
class PipelineStatistics
```

A set of statistics for the processed document

documents\_processed document\_exceptions

```
processed_document (document)
```

Update statistics based on this document completing processing

**Parameters** **document** – the document that has been processed

## 3.5 Sinks

Sinks are the end-point of a Pipeline and allow for the final output of the pipeline to be either stored or written out

```
class InMemoryDocumentSink
```

An in-memory document sink can be used for testing where you want to capture a set of the documents as basic list in-memory and then access them

```
get_document (index)
```

Get document at given index

**Parameters** **index** – index to get the document at

**sink** (*document*)

Adds the document to the sink

**Parameters** **document** – document to add

## 3.6 Stores

Stores are persistence components for Documents, typically they can act as either a Connector or a Sink

**class** **JsonDocumentStore** (*store\_path: str, force\_initialize: bool = False*)

An implementation of a document store that uses JSON files to store the documents and maintains an index.json containing some basics of the documents

**add** (*document: kodexa.model.model.Document*)

Add a new document and return the index position

**Returns** The index of the document added

**count** ()

The number of documents in the store

**Returns** The number of documents

**delete** (*idx: int*)

Delete the document at the given index

**Returns** The Document that was removed

**get** (*idx: int*)

Load the document at the given index

**Returns** Document at given index

**get\_document** (*index: int*)

Gets the document from the specific index

**Parameters** **index** – index of document to get

**Returns** the document

**load** (*document\_id: str*)

Loads the document with the given document ID

:return the document

**read\_index** ()

Method to read the document index from the store path

**reset\_connector** ()

Reset the index back to the beginning

**save\_index** ()

Method to write the JSON store index back to store path

**class** **TableDataStore** (*columns=None, rows=None*)

Stores data as a list of lists that can represent a table.

This is a good store when you are capturing nested or tabular data.

**Parameters**

- **columns** – a list of the column names (default to dynamic)



- **rows** – initial set of rows (default to empty)

**add** (*row*)

Writes a row to the Data Store

**Parameters** **row** – the row (as a list) to add

**count** ()

Returns the number of rows in the store

**Returns** number of rows

**class DictDataStore** (*dicts=None*)

Stores data as a list of dictionaries that can be any structure

This is a good store when you are capturing nested or semi-structured data

**add** (*dict*)

Writes a dict to the Data Store

**Parameters** **dict** – the dict to add to the store

**count** ()

Returns the number of dictionaries in the store

**Returns** number of dictionaries

**class NodeTagger** (*type\_re, content\_re, tag\_name, use\_all\_content=True, node\_only=False*)

A node tagger allows you to provide a type and content regular expression and then tag content in all matching nodes.

It allows for multiple matching groups to be defined, also the ability to use all content and also just tag the node (ignoring the matching groups)

**class TextParser** (*decode=False, encoding='utf-8'*)

The text parser can load a source file as a text document and creates a single content node with the text

**class Rollup** (*collapse\_type\_res=[], reindex=True*)

The rollup step allows you to decide how you want to collapse content in a document by removing nodes while maintaining content and features as needed

**class ExtractTagsToKeyValuePair** (*store\_name, include=[], exclude=[], include\_node\_content=True*)

Extract all the tags from a document into a key/value pair table store

**class JsonParser**

Parse JSON file into kodexa Document

**class ExtractTagsToKeyValuePair** (*store\_name, include=[], exclude=[], include\_node\_content=True*)

Extract all the tags from a document into a key/value pair table store

**class JsonParser**

Parse JSON file into kodexa Document

**class NodeTagger** (*type\_re, content\_re, tag\_name, use\_all\_content=True, node\_only=False*)

A node tagger allows you to provide a type and content regular expression and then tag content in all matching nodes.

It allows for multiple matching groups to be defined, also the ability to use all content and also just tag the node (ignoring the matching groups)

**class Rollup** (*collapse\_type\_res=[], reindex=True*)

The rollup step allows you to decide how you want to collapse content in a document by removing nodes while maintaining content and features as needed

**class TextParser** (*decode=False, encoding='utf-8'*)

The text parser can load a source file as a text document and creates a single content node with the text

## CHAPTER 4

---

### Indices and tables

---

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



**k**

`kodexa`, 13  
`kodexa.cloud`, 13  
`kodexa.connectors`, 13  
`kodexa.model`, 13  
`kodexa.pipeline`, 17  
`kodexa.sinks`, 19  
`kodexa.steps.common`, 21  
`kodexa.stores`, 20



**A**

add() (*DictDataStore* method), 21  
 add() (*JsonDocumentStore* method), 20  
 add() (*TableDataStore* method), 21  
 add\_child() (*ContentNode* method), 15  
 add\_feature() (*ContentNode* method), 15  
 add\_mixin() (*Document* method), 14  
 add\_step() (*Pipeline* method), 18  
 add\_store() (*Pipeline* method), 18  
 add\_store() (*PipelineContext* method), 19  
 adopt\_children() (*in module kodexa.mixins.core*),  
 8

**C**

collapse() (*in module kodexa.mixins.spatial*), 12  
 collect\_nodes\_to() (*in module kodexa.mixins.core*), 9  
 ContentNode (*class in kodexa.model*), 15  
 count() (*DictDataStore* method), 21  
 count() (*JsonDocumentStore* method), 20  
 count() (*TableDataStore* method), 21  
 create\_node() (*Document* method), 14

**D**

delete() (*JsonDocumentStore* method), 20  
 DictDataStore (*class in kodexa.stores*), 21  
 Document (*class in kodexa.model*), 14  
 DocumentMetadata (*class in kodexa.model*), 15

**E**

ExtractTagsToKeyValuePair (*class in kodexa.steps.common*), 21

**F**

find() (*in module kodexa.mixins.core*), 6  
 find\_with\_feature\_value() (*in module kodexa.mixins.core*), 7  
 findall() (*in module kodexa.mixins.core*), 7

findall\_with\_feature\_value() (*in module kodexa.mixins.core*), 7

FolderConnector (*class in kodexa.connectors*), 13  
 from\_dict() (*Document* static method), 14  
 from\_json() (*Document* static method), 14  
 from\_kdxa() (*Document* static method), 14  
 from\_msgpack() (*Document* static method), 14

**G**

get() (*JsonDocumentStore* method), 20  
 get\_all\_content() (*in module kodexa.mixins.core*), 6  
 get\_all\_tags() (*in module kodexa.mixins.core*), 6  
 get\_bbox() (*in module kodexa.mixins.spatial*), 11  
 get\_children() (*ContentNode* method), 16  
 get\_content() (*ContentNode* method), 16  
 get\_document() (*InMemoryDocumentSink* method),  
 19  
 get\_document() (*JsonDocumentStore* method), 20  
 get\_feature() (*ContentNode* method), 16  
 get\_feature\_value() (*ContentNode* method), 16  
 get\_features() (*ContentNode* method), 16  
 get\_features\_of\_type() (*ContentNode* method),  
 16  
 get\_height() (*in module kodexa.mixins.spatial*), 12  
 get\_last\_child\_index() (*in module kodexa.mixins.core*), 10  
 get\_mixins() (*Document* method), 14  
 get\_node\_at\_index() (*in module kodexa.mixins.core*), 9  
 get\_root() (*Document* method), 14  
 get\_rotate() (*in module kodexa.mixins.spatial*), 11  
 get\_statistics() (*in module kodexa.mixins.spatial*), 11  
 get\_store() (*PipelineContext* method), 19  
 get\_store\_names() (*PipelineContext* method), 19  
 get\_tag() (*in module kodexa.mixins.core*), 6  
 get\_tags() (*in module kodexa.mixins.core*), 8  
 get\_type() (*ContentNode* method), 16  
 get\_width() (*in module kodexa.mixins.spatial*), 12

get\_x() (in module *kodexa.mixins.spatial*), 11  
get\_y() (in module *kodexa.mixins.spatial*), 11

## H

has\_feature() (*ContentNode* method), 16  
has\_next\_node() (in module *kodexa.mixins.core*), 9  
has\_previous\_node() (in module *kodexa.mixins.core*), 9

## I

InMemoryDocumentSink (class in *kodexa.sinks*), 19  
is\_first\_child() (in module *kodexa.mixins.core*), 10  
is\_last\_child() (in module *kodexa.mixins.core*), 10

## J

JsonDocumentStore (class in *kodexa.stores*), 20  
JsonParser (class in *kodexa.steps.common*), 21

## K

kodexa (module), 13  
kodexa.cloud (module), 13  
kodexa.connectors (module), 13  
kodexa.model (module), 13  
kodexa.pipeline (module), 17  
kodexa.sinks (module), 19  
kodexa.steps.common (module), 21  
kodexa.stores (module), 20  
KodexaCloudPipeline (class in *kodexa.cloud*), 13  
KodexaCloudService (class in *kodexa.cloud*), 13

## L

load() (*JsonDocumentStore* method), 20

## M

move\_child\_to\_parent() (in module *kodexa.mixins.core*), 8

## N

next\_node() (in module *kodexa.mixins.core*), 9  
NodeTagger (class in *kodexa.steps.common*), 21

## P

Pipeline (class in *kodexa.pipeline*), 17  
PipelineContext (class in *kodexa.pipeline*), 19  
PipelineStatistics (class in *kodexa.pipeline*), 19  
previous\_node() (in module *kodexa.mixins.core*), 10  
processed\_document() (*PipelineStatistics* method), 19

## R

read\_index() (*JsonDocumentStore* method), 20

remove\_feature() (*ContentNode* method), 17  
remove\_tag() (in module *kodexa.mixins.core*), 8  
reset\_connector() (*JsonDocumentStore* method), 20

Rollup (class in *kodexa.steps.common*), 21  
run() (*Pipeline* method), 18

## S

save\_index() (*JsonDocumentStore* method), 20  
set\_bbox() (in module *kodexa.mixins.spatial*), 11  
set\_bbox\_from\_children() (in module *kodexa.mixins.spatial*), 12  
set\_feature() (*ContentNode* method), 17  
set\_output\_document() (*PipelineContext* method), 19  
set\_rotate() (in module *kodexa.mixins.spatial*), 11  
set\_sink() (*Pipeline* method), 18  
set\_statistics() (in module *kodexa.mixins.spatial*), 10  
sink() (*InMemoryDocumentSink* method), 20

## T

TableDataStore (class in *kodexa.stores*), 20  
tag() (in module *kodexa.mixins.core*), 5  
tag\_nodes\_to() (in module *kodexa.mixins.core*), 9  
tag\_range() (in module *kodexa.mixins.core*), 5  
TextParser (class in *kodexa.steps.common*), 21  
to\_dict() (*ContentNode* method), 17  
to\_dict() (*Document* method), 14  
to\_json() (*ContentNode* method), 17  
to\_json() (*Document* method), 15  
to\_kdxa() (*Document* method), 15  
to\_msgpack() (*Document* method), 15  
to\_text() (*ContentNode* method), 17  
to\_text() (*Document* method), 15

## U

UrlConnector (class in *kodexa.connectors*), 13