# intake_avro Documentation

***Release 0.0.1***

**Joseph Crail**

**Dec 03, 2018**

# Contents:

This package enables the loading of Apache Avro files within the Intake data loading and catalog system. Two plugins are contained: for fast loading of strictly tabular data to pandas dataframes, and slower reading of more complicatedly structured data as a sequence of python dictionaries.

Each avro file becomes one partition.

# Quickstart

`intake_avro` provides quick and easy access to tabular data stored in the Apache Avro binary, columnar format.

## 1.1 Installation

To use this plugin for intake, install with the following command:

```
conda install -c intake intake-avro
```

## 1.2 Usage

### 1.2.1 Ad-hoc

After installation, the functions `intake.open_avro_table` and `intake.open_avro_sequence` will become available. The former, faster method can be used to open one or more Avro files with *flat* schema into dataframes, but the latter can be used for any files and produces generic sequences of dictionaries.

Assuming some Avro files in a given path, the following would load them into a dataframe:

```python
import intake
source = intake.open_avro_table('data_path/*.avro')
dataframe = source.read()
```

There will be one data partition per input file; there is no random access within each Avro data file.

Arguments to the `open_avro_*` functions:

- `urlpath` : the location of the data. This can be a single file, a list of specific files,

    or a glob string (containing `"*"`). The URLs can be local files or, if using a protocol specifier such as `'s3://'`, a remote file location.

---

- `storage_options` : other parameters that are to be passed to the filesystem

  implementation, in the case that a remote filesystem is referenced in `urlpath`. For specifics, see the Dask documentation.

A source so defined will provide the usual methods such as `discover` and `read_partition`.

### 1.2.2 Creating Catalog Entries

To include in a catalog, the plugin must be listed in the plugins of the catalog:

```
plugins:
  source:
    - module: intake_avro
```

and entries must specify `driver: avro_table` or `driver: avro_sequence`. The further arguments are exactly the same as for the `open_avro_*` functions.

### 1.2.3 Using a Catalog

Assuming a catalog file called `cat.yaml`, containing a Avro source `pdata`, one could load it into a dataframe as follows:

```python
import intake
cat = intake.Catalog('cat.yaml')
df = cat.pdata.read()
```

The type of the output will depend on the plugin that was defined in the catalog. You can inspect this before loading by looking at the `.container` attribute, which will be either `"dataframe"` or `"python"`.

The number of partitions will be equal to the number of files pointed to.

# API Reference

| | |
|---|---|
| *intake_avro.source.AvroTableSource*(urlpath) | Source to load tabular Avro datasets. |
| *intake_avro.source.AvroSequenceSource*(urlpath) | Source to load Avro datasets as sequence of Python dicts. |

**class** intake_avro.source.**AvroTableSource**(*urlpath*, *metadata=None*, *storage_options=None*)

Source to load tabular Avro datasets.

> **Parameters**
>
> > **urlpath: str** Location of the data files; can include protocol and glob characters.
>
> **Attributes**
>
> > **cache_dirs**
> >
> > **datashape**
> >
> > **description**
> >
> > **hvplot** Returns a hvPlot object to provide a high-level plotting API.
> >
> > **plot** Returns a hvPlot object to provide a high-level plotting API.
> >
> > **plots** List custom associated quick-plots

### Methods

| | |
|---|---|
| close() | Close open resources corresponding to this data source. |
| discover() | Open resource and populate the source attributes. |
| *read*() | Load entire dataset into a container and return it |

Table 2 – continued from previous page

| read_chunked() | Return iterator over container fragments of data source |
|---|---|
| read_partition(i) | Return a (offset_tuple, container) corresponding to i-th partition. |
| *to_dask*() | Create lazy dask dataframe object |
| *to_spark*() | Pass URL to spark to load as a DataFrame |
| yaml([with_plugin]) | Return YAML representation of this data-source |

set_cache_dir

**read**()
  Load entire dataset into a container and return it

**to_dask**()
  Create lazy dask dataframe object

**to_spark**()
  Pass URL to spark to load as a DataFrame

  Note that this requires org.apache.spark.sql.avro.AvroFileFormat to be installed in your spark classes.

  This feature is experimental.

**class** intake_avro.source.**AvroSequenceSource**(*urlpath*, *metadata=None*, *storage_options=None*)
  Source to load Avro datasets as sequence of Python dicts.

  **Parameters**

  **urlpath: str** Location of the data files; can include protocol and glob characters.

  **Attributes**

  **cache_dirs**

  **datashape**

  **description**

  **hvplot** Returns a hvPlot object to provide a high-level plotting API.

  **plot** Returns a hvPlot object to provide a high-level plotting API.

  **plots** List custom associated quick-plots

  ### Methods

| close() | Close open resources corresponding to this data source. |
|---|---|
| discover() | Open resource and populate the source attributes. |
| *read*() | Load entire dataset into a container and return it |
| read_chunked() | Return iterator over container fragments of data source |
| read_partition(i) | Return a (offset_tuple, container) corresponding to i-th partition. |

Table 3 – continued from previous page

| | |
|---|---|
| *to_dask*() | Create lazy dask bag object |
| to_spark() | Provide an equivalent data object in Apache Spark |
| yaml([with_plugin]) | Return YAML representation of this data-source |

| set_cache_dir | |
|---|---|

**read**()
  Load entire dataset into a container and return it

**to_dask**()
  Create lazy dask bag object

# Indices and tables

- genindex
- modindex
- search

# Index

## A

AvroSequenceSource (*class in intake_avro.source*), [6]

AvroTableSource (*class in intake_avro.source*), [5]

## R

read() (*intake_avro.source.AvroSequenceSource method*), [7]

read() (*intake_avro.source.AvroTableSource method*), [6]

## T

to_dask() (*intake_avro.source.AvroSequenceSource method*), [7]

to_dask() (*intake_avro.source.AvroTableSource method*), [6]

to_spark() (*intake_avro.source.AvroTableSource method*), [6]