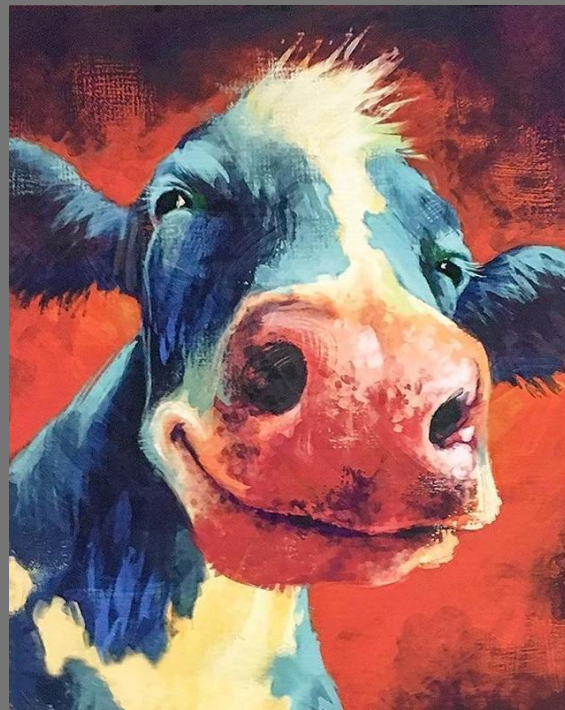

Cheap Long-Term Storage of structured data

Using AWS Data Wrangler, S3 and Parquet files

@enricomarchesin
Principal Software Architect
@Cainthus



User Requirements

- structured data
 - long-term storage
 - large* batch processing oriented
 - standard & Open*
 - cheap
-

Ideas: data format

- structured data
- long-term storage
- large* batch processing oriented
- standard & Open*
- cheap

AVRO or PARQUET

Ideas: storage

- structured data
- long-term storage
- large* batch processing oriented
- standard & Open*
- cheap

The logo for Amazon S3, consisting of the letters 'S' and '3' in a stylized, rounded, orange font.

Tools

- AVRO/Parquet
 - structured data
 - large* batch processing oriented
 - standard & Open*
 - S3
 - long-term storage
 - standard & Open*
 - cheap
-

Tools (in Python)

- AVRO/Parquet
 - structured data
 - large* batch processing oriented
 - standard & Open*

- S3

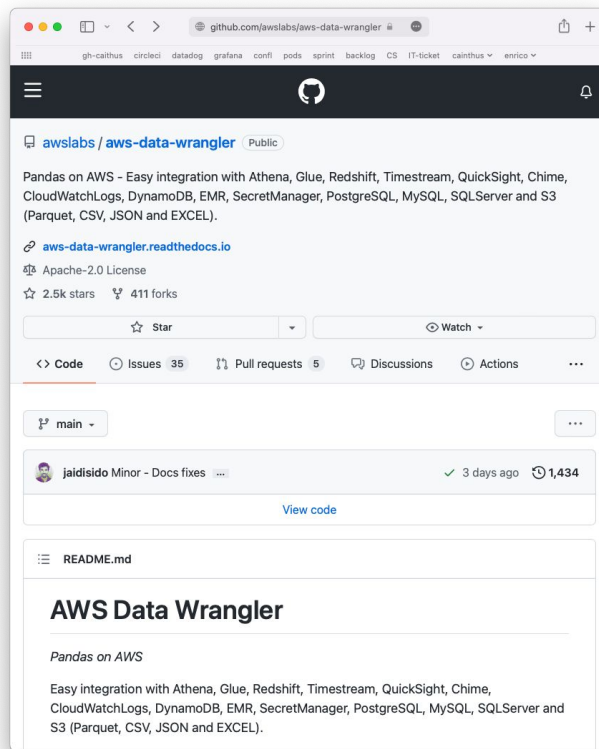
- long-term storage
- standard & Open
- cheap

Pandas

Tech requirements

- AVRO/Parquet
- S3
- the API is a Pandas DataFrame
- standard and Open*

AWS Data Wrangler



<https://github.com/aws-labs/aws-data-wrangler>

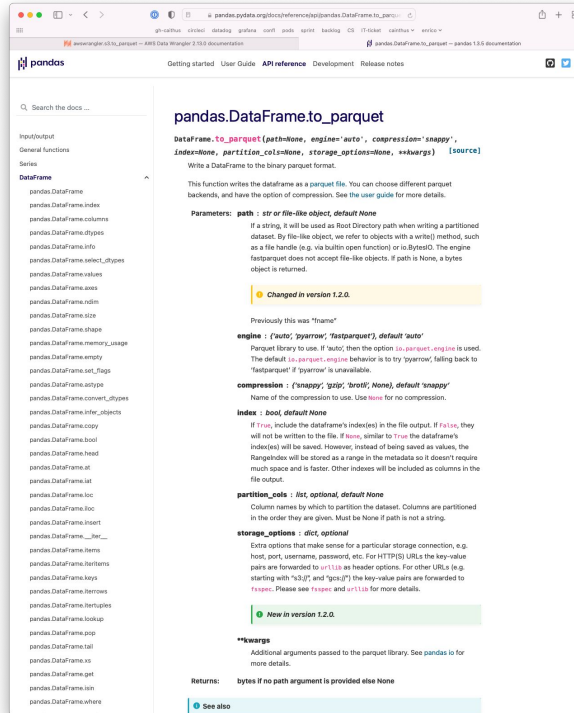


Let's code!

Tenets Check

- AWS Data Wrangler
 - AVRO/Parquet
 - S3
 - the API is a Pandas DataFrame
 - standard and Open*

From Pandas...



...to ADW

awsdatawrangler.s3.to_parquet

The concept of Dataset goes beyond the simple idea of ordinary files and enable more complex features like partitioning and catalog integration (Amazon Athena/HMS Glue Catalog).

This operation may mutate the original pandas dataframe in place. To avoid this behaviour please pass in a deep copy instead (i.e. `df.copy()`).

Database and table arguments are optional, the table name and all column names will be automatically sanitized using `wr.catalog.sanitize_table_name` and `wr.catalog.sanitize_column_name`. Please, see `sanitize_column_name` to enforce this behaviour deep.

In append mode, the parameters will be spent on an existing table.

In case of `use_threads=True` the number of threads that will be spawned will be gotten from `os.cpu_count()`.

This function has arguments which can be configured globally through `wr.config` or environment variables:

- `catalog_id`
- `concurrent_partitioning`
- `dataset`

Check out the [Global Configuration Tutorial](#) for details.

Parameters

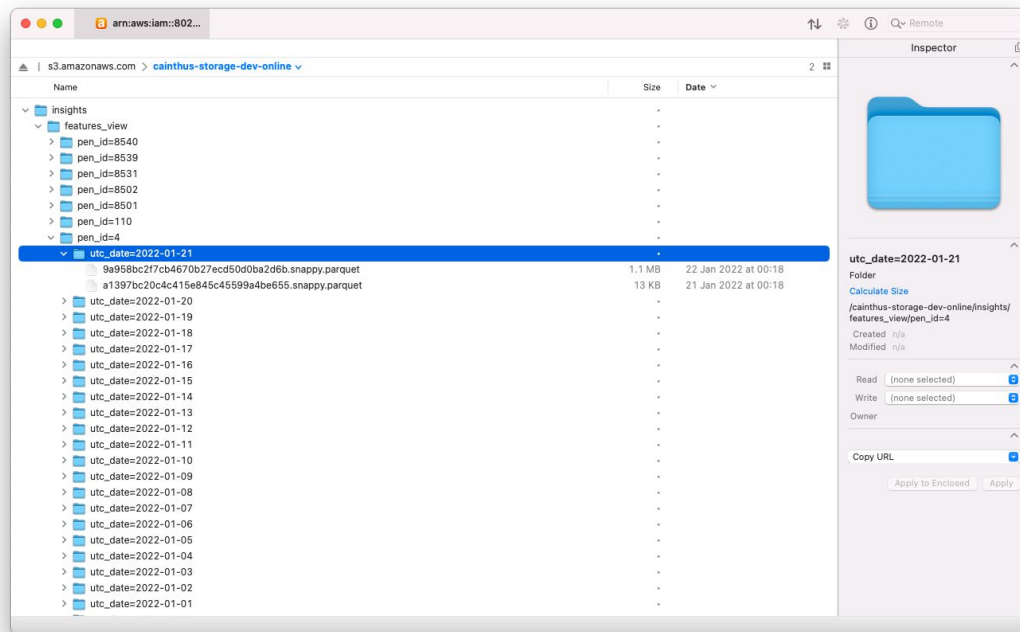
- df** (pandas.DataFrame) – Pandas Dataframe <https://pandas.pydata.org/pandas-docs/stable/references/api/pandas.DataFrame.html>
- path** (str, optional) – S3 path like file e.g. `s3://bucket/path/to/s3/basename.parquet` for dataset e.g. `s3://bucket/path/1`. Required if dataset=False or when dataset=True and creating a new dataset
- index** (bool) – True to store the DataFrame index as file, otherwise False to ignore it.
- compression** (str, optional) – Compression type (None, GZIP, BZIP2, SNAPPY, ZSTD).
- pyarrow_additional_kwargs** (Optional[Dict[str, Any]]) – Additional parameters forwarded to pyarrow e.g. `pyarrow_additional_kwargs={"coerce_timestamps": "us", "use_deprecated_int64_timestamps": False, "allow_increased_timestamp_precision": False}`
- max_rows_by_file** (int) – Max number of rows in each file. Default is None i.e. don't split the files. In e.g. 33334442, 344434454
- use_threads** (bool, int) – True to enable concurrent requests, False to disable multiple threads. Forbidding `use_threads` will be used as the max number of threads. If integer is provided, specified number is used.
- bucket_location** (Optional[Dict[str, Any]]) – Bucket Location. The default bucket location will be used if `bucket_location` is None.
- s3_additional_kwargs** (Optional[Dict[str, Any]]) – Forwarded to boto3 requests, e.g. `s3_additional_kwargs={"ServerSideEncryption": "aws:kms", "SSVtoSSEKey": "TODOR_KMS_KEY_ABT"}`
- sanitize_column_names** (bool) – True to sanitize column names (using `wr.catalog.sanitize_table_name` and `wr.catalog.sanitize_column_name`) or False to keep it as is. True value behaviour is enforced if database and table arguments are passed.
- dataset** (bool) – If True store a parquet dataset instead of a ordinary file(s). True, enable all follow arguments: `partition_cols`, `mode`, `database`, `table`, `description`, `parameters`, `column_comments`, `concurrent_partitioning`, `catalog`, `sanitizing`, `projection_enabled`, `projection_type`, `projection_ranges`, `projection_values`, `projection_intervals`, `concurrent_partitioning`, `catalog`, `sanitizing`, `projection_enabled`, `projection_type`, `projection_ranges`, `projection_values`, `projection_intervals`.
- database** (str, optional) – If dataset=True, add a database prefix to the output files.
- partition_cols** (List[str], optional) – List of column names that will be used to create partitions. Only takes effect if dataset=True.
- bucketing_info** (Tuple[str], optional) – Tuple consisting of the column names used for bucketing as the first element and the number of buckets as the second element. Only str, int and bool are supported as column data types for bucketing.
- concurrent_partitioning** (bool) – If True will increase the partitioning level during the partitioning writing. It will decrease the writing time and increase the memory usage. <https://aws-data-wrangler.readthedocs.io/en/2.13.0/tutorials/22%20-%20Accelerating%20Partitioning%20Concurrently.html>
- make_id** (str, optional) – Internal ID used by the function. Only takes effect if dataset=True. For details check the related tutorial <https://aws-data-wrangler.readthedocs.io/en/2.13.0/tutorials/22%20-%20Accelerating%20Partitioning%20Concurrently.html>
- catalog** (bool) – If True and mode="overwrite", creates an archived version of the table catalog before updating it.
- schema_evolution** (bool) – If True, allows schema evolution from a source column to a different one, and it will be added. True by default. Only considered

From `pd` to `wr`

```
df.to_parquet(  
    path=f"s3://my-bucket/my-path/",  
    partition_cols=["date", "obj_id"],  
)
```

```
wr.s3.to_parquet(  
    df,  
    path=f"s3://my-bucket/my-path/",  
    partition_cols=["date", "obj_id"],  
    dataset=True,  
)
```

Partitions



AWS Glue?

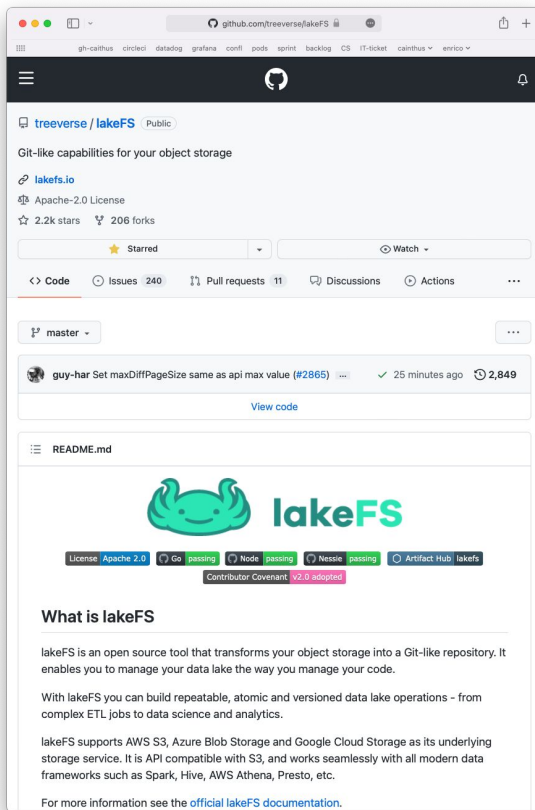
Maybe

But we don't need it now

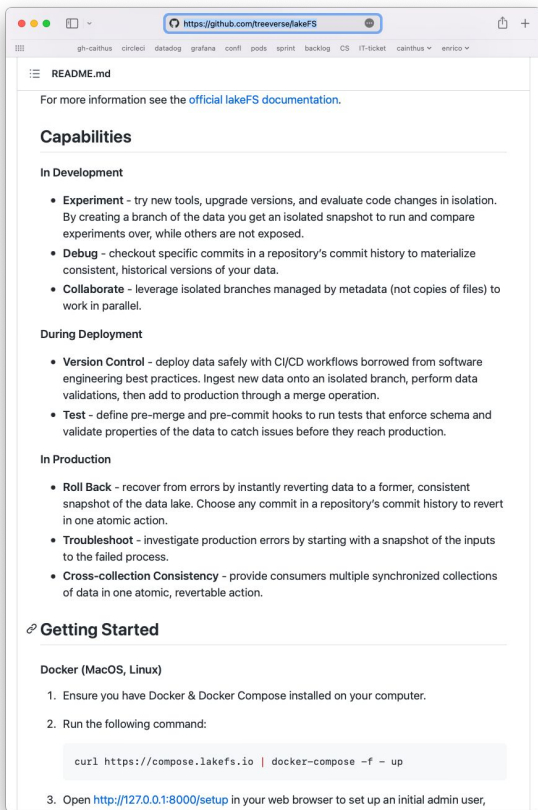
KISS 🙄



What's next...



<https://github.com/treeverse/lakeFS>



https://github.com/treeverse/lakeFS