

---

# Feast Spark SDK

**Feast Authors**

**Jun 21, 2022**



**CONTENTS:**

<b>1 Client</b>	<b>1</b>
<b>Python Module Index</b>	<b>5</b>
<b>Index</b>	<b>7</b>



```
class feast_spark.client.Client(feast_client: Client)
```

Bases: object

```
property config: Config
```

```
property feature_store: Client
```

```
get_health_metrics(project: str, table_names: List[str]) → Dict[str, List[str]]
```

```
get_historical_features(feature_refs: List[str], entity_source: Union[DataFrame, FileSource, BigQuerySource], output_location: Optional[str] = None) → RetrievalJob
```

Launch a historical feature retrieval job.

#### Parameters

- **feature\_refs** – List of feature references that will be returned for each entity. Each feature reference should have the following format: “feature\_table:feature” where “feature\_table” & “feature” refer to the feature and feature table names respectively.
- **entity\_source** (*Union[`pd.DataFrame`, `FileSource`, `BigQuerySource`]*) – Source for the entity rows. If `entity_source` is a `Panda DataFrame`, the dataframe will be staged to become accessible by spark workers. If one of feature tables’ source is in `BigQuery` - entities will be upload to `BQ`. Otherwise to remote file storage (derived from configured staging location). It is also assumed that the column `event_timestamp` is present in the dataframe, and is of type `datetime` without `timezone` information.

The user needs to make sure that the source (or staging location, if `entity_source` is a `Panda DataFrame`) is accessible from the Spark cluster that will be used for the retrieval job.

- **output\_location** – Specifies the path in a bucket to write the exported feature data files

#### Returns

Returns a retrieval job object that can be used to monitor retrieval progress asynchronously, and can be used to materialize the results.

## Examples

```
>>> import feast
>>> import feast_spark
>>> from feast.data_format import ParquetFormat
>>> from datetime import datetime
>>> feast_client = feast.Client(core_url="localhost:6565")
>>> feature_refs = ["bookings:bookings_7d", "bookings:booking_14d"]
>>> entity_source = FileSource("event_timestamp", ParquetFormat(), "gs://some-
↳bucket/customer")
>>> feature_retrieval_job = feast_spark.Client(feast_client).get_historical_
↳features(
>>>     feature_refs, entity_source)
>>> output_file_uri = feature_retrieval_job.get_output_file_uri()
"gs://some-bucket/output/"
```

`get_historical_features_df`(*feature\_refs: List[str]*, *entity\_source: Union[FileSource, BigQuerySource]*)

Launch a historical feature retrieval job.

### Parameters

- **feature\_refs** – List of feature references that will be returned for each entity. Each feature reference should have the following format: “feature\_table:feature” where “feature\_table” & “feature” refer to the feature and feature table names respectively.
- **entity\_source** (*Union[FileSource, BigQuerySource]*) – Source for the entity rows. The user needs to make sure that the source is accessible from the Spark cluster that will be used for the retrieval job.

### Returns

Returns the historical feature retrieval result in the form of Spark dataframe.

## Examples

```
>>> import feast
>>> import feast_spark
>>> from feast.data_format import ParquetFormat
>>> from datetime import datetime
>>> from pyspark.sql import SparkSession
>>> spark = SparkSession.builder.getOrCreate()
>>> feast_client = feast.Client(core_url="localhost:6565")
>>> feature_refs = ["bookings:bookings_7d", "bookings:booking_14d"]
>>> entity_source = FileSource("event_timestamp", ParquetFormat(), "gs://some-
↳bucket/customer")
>>> df = feast_spark.Client(feast_client).get_historical_features(
>>>     feature_refs, entity_source)
```

`get_job_by_id`(*job\_id: str*) → SparkJob

`list_jobs`(*include\_terminated: bool*, *project: Optional[str] = None*, *table\_name: Optional[str] = None*) → List[SparkJob]

List ingestion jobs currently running in Feast.

### Parameters

- **include\_terminated** – Flag to include terminated jobs or not
- **project** – Optionally specify the project to use as filter when retrieving jobs
- **table\_name** – Optionally specify name of feature table to use as filter when retrieving jobs

**Returns**

List of SparkJob ingestion jobs.

**property metrics\_redis: Redis**

**schedule\_offline\_to\_online\_ingestion**(*feature\_table: FeatureTable, ingestion\_timespan: int, cron\_schedule: str*)

Launch Scheduled Ingestion Job from Batch Source to Online Store for given feature table

**Parameters**

- **feature\_table** – FeatureTable that will be ingested into the online store
- **ingestion\_timespan** – Days of data which will be ingestion per job. The boundaries on which to filter the source are [end of day of execution date - ingestion\_timespan (days) , end of day of execution date)
- **cron\_schedule** – Cron schedule expression

Returns: Spark Job Proxy object

**start\_offline\_to\_online\_ingestion**(*feature\_table: FeatureTable, start: datetime, end: datetime*) → SparkJob

Launch Ingestion Job from Batch Source to Online Store for given feature table

**Parameters**

- **feature\_table** – FeatureTable that will be ingested into the online store
- **start** – lower datetime boundary on which to filter the source
- **end** – upper datetime boundary on which to filter the source

Returns: Spark Job Proxy object

**start\_stream\_to\_online\_ingestion**(*feature\_table: FeatureTable, extra\_jars: Optional[List[str]] = None, project: Optional[str] = None*) → SparkJob

**unschedule\_offline\_to\_online\_ingestion**(*feature\_table: FeatureTable, project=None*)

**feast\_spark.client.stage\_entities\_to\_bq\_with\_partition**(*entity\_source: DataFrame, project: str, dataset: str*) → BigQuerySource

Stores given (entity) dataframe as new table in BQ. Name of the table generated based on current time. Table will expire in 1 day. Returns BigQuerySource with reference to created table.





## PYTHON MODULE INDEX

f

`feast_spark.client`, 1



## C

`Client` (*class in feast\_spark.client*), 1  
`config` (*feast\_spark.client.Client property*), 1

## F

`feast_spark.client`  
    module, 1  
`feature_store` (*feast\_spark.client.Client property*), 1

## G

`get_health_metrics()` (*feast\_spark.client.Client method*), 1  
`get_historical_features()`  
    (*feast\_spark.client.Client method*), 1  
`get_historical_features_df()`  
    (*feast\_spark.client.Client method*), 2  
`get_job_by_id()` (*feast\_spark.client.Client method*), 2

## L

`list_jobs()` (*feast\_spark.client.Client method*), 2

## M

`metrics_redis` (*feast\_spark.client.Client property*), 3  
module  
    `feast_spark.client`, 1

## S

`schedule_offline_to_online_ingestion()`  
    (*feast\_spark.client.Client method*), 3  
`stage_entities_to_bq_with_partition()` (*in module feast\_spark.client*), 3  
`start_offline_to_online_ingestion()`  
    (*feast\_spark.client.Client method*), 3  
`start_stream_to_online_ingestion()`  
    (*feast\_spark.client.Client method*), 3

## U

`unschedule_offline_to_online_ingestion()`  
    (*feast\_spark.client.Client method*), 3