**Introduction of AOTAS  modules and source code**

● Module Overview

| Module Name | Description |
|---|---|
| Download | Download data files To download dataset data files from source and save  to **'raw'** folder, using the original format, **'csv'.** |
| ETL Module | Read 'raw' data, perform cleansing transformation and save to 'processed' location |
| Integration Module | Read 'processed' data,  conduct data warehouse modeling, create dimension/fact tables and save to 'integration' location |
| Utility Module | Provide common utilities that will be used throughout the pipeline process |

● Module details

| Module | Scripts | | |
|---|---|---|---|
| Top Entry point | .AOTA_main.py.<br>● Entry point of the pipeline<br>● It  will run the main  file of download module,   ETL module , and integration module to complete their duties. | | |
| Download Module | Description | To download dataset data files from source and save  to **'raw'** folder, using the original format, **'csv'.** | |
| | Entry point file | **download_manager.py** | |
| | Major Files | **downloader.p**y– to  provide tools/methods  to download flight time dataset<br><br>Two major methods | |
| | | Method | Description |

| | | download_small_dataset | To download following small datasets: airport, carrier, plane |
|---|---|---|---|
| | | download_dataset_flight | To download flight time dataset of multiple years (specified in array `flight_dfiles`) |

| ETL Module | Description | Provide following functionalities:<br>1. Read data from **'raw'** files,<br>2. perform cleansing/transformation<br>3. save to **'processed'** folder, using '**parquet**' |
|---|---|---|
| | EntryPoint File | **ETL_Manager.py** |
| | Major Files | **ETL_Airport.py**: to process airport data<br>**ETL_Carrier.py:** to process carrier data<br>**ETL_Plane.py**: to process plane data<br>**ETL_Flight.py**:to process flight data |

| Integration Module | Description | Provide following functionalities :<br>1. Read data from **'process'** files,<br>2. Perform DW modeling,<br>3. Create dimension table and fact table datasets<br>4. Save to **'integration'** folder, using '**parquet**' |
|---|---|---|
| | EntryPoint File | **intgr_Manager.py** |
| | Major Files | **1.Intgr_stg_table.py** (main method: `cr_stg_flight`)<br>Read 'processed' flight data and create a staging table dataset, **stg_flight**, which works as a temp table for downstream processes. It will be saved to "integrated" folder, usin "parquet"<br><br>**2.Intgr_dim_table.py**<br> Read data from "processed' folder and use 'STG_FLIGHT' data created in step 1, create following dimension table datasets and save to "integration"folder, using "parquet".<br>**dim_plane**<br>**dim_origin_dest**<br>**dim_date**<br>**dim_carrier** |

| | | |
|---|---|---|
| | | **3.Intgr_fact_table.py**(main method: **cr_fact_flight**)<br>    1) Read from **stg_flight** (created in **step 1**) and **dim_\* datasets** (Created in **step 2**),<br>    2) Create the  flight fact table dataset **'fact_flight'**<br>    3) Save to  "integration"folder, using "parquet". |
| Utility module | Description | Provide different types of  utilities to help the whole process. |
| | Major files | **ut_store.py:**  work as a configuration file, which stores various parameters (file name,process name,etc).<br><br>**ut_base.py**: provide some most  basic utilities that will used by all modules<br><br>**ut_log.py**:  provide logging related functionalities that will be used by Download, ETL and interaction modules<br><br>**Ut_spark.py:** Provided utilities that are related to Spark, such as, create a spark session, read from/write to files, etc<br><br>**Ut_pipeline.py:** Provided utilities that are related to pipeline processing, such as log metrics, read data from 'raw', or save to 'processed', etc. They will be mainly  used by ETL and Integration modules. |