

< [Return to "Data Engineering Nanodegree" in the classroom](#)

Data Pipelines with Airflow

REVIEW	CODE REVIEW	HISTORY
<h2>Meets Specifications</h2> <p>IMPRESSIVE!! You did very well on this and passed the project successfully from the first trail! Congratulations :) Thanks a lot for your effort and commitment to get the project completed. wish you the best of luck in future submissions and projects :)</p>		
<h3>General</h3>		
<div><div>✓</div><div>DAG can be browsed without issues in the Airflow UI</div></div>		
Correct! Dag loaded and Browsed successfully through the interface without any issues		
<div><div>✓</div><div>The dag follows the data flow provided in the instructions, all the tasks have a dependency and DAG begins with a start_execution task and ends with a end_execution task.</div></div>		
Good work! DAG follows the required data flow provided. All functions have the correct dependencies, star point and end point are present as requested in the instructions		
<h3>Dag configuration</h3>		
<div><div>✓</div><div>DAG contains default_args dict, with the following keys:<ul style="list-style-type: none">OwnerDepends_on_pastStart_dateRetriesRetry_delayCatchup</div></div>		
All default arguments are defined correctly as required in the default_args dict		
<div><div>✓</div><div>The DAG object has default args set</div></div>		
DAG defined in the right manner, and the default_args dict was assigned in the right manner to the DAG's default_args parameter		
<div><div>✓</div><div>The DAG should be scheduled to run once an hour</div></div>		
Excellent! scheduling was done using <code>0 * * * * *</code> . You can also use the <code>@hourly</code> expression		
<h3>Staging the data</h3>		
<div><div>✓</div><div>There is a task that to stages data from S3 to Redshift. (Runs a Redshift copy statement)</div></div>		
Very good! the StageToRedshiftOperator is present as required and it's running the copy statements from S3 o AWS Redshift		
<div><div>✓</div><div>Instead of running a static SQL statement to stage the data, the task uses params to generate the copy statement dynamically</div></div>		
Great work! The copy statement is parametrized and not hardcoded, so it can load multiple tables with different config and not tied to only table.		
<div><div>✓</div><div>The operator contains logging in different steps of the execution</div></div>		
Good work but only two lines of logging is present. Logging should be enhanced to be able to monitor the different functionalities occurring within your operator and to be able to debug your code in case of any issue as this will minimize the time required to resolve an issue if any.		
<div><div>✓</div><div>The SQL statements are executed by using a Airflow hook</div></div>		
Nice work on this. Used the PostgresHook to define your connection		
<h3>Loading dimensions and facts</h3>		
<div><div>✓</div><div>Dimensions are loaded with on the LoadDimension operator</div></div>		
Nice work! Dimensions are loaded using the LoadDimensionOperator operator successfully		
<div><div>✓</div><div>Facts are loaded with on the LoadFact operator</div></div>		
Correct! You used the loadFactOperator operator in the correct manner required.		
<div><div>✓</div><div>Instead of running a static SQL statement to stage the data, the task uses params to generate the copy statement dynamically</div></div>		
You did well on this by passing SQL queries as a parameter with the table and connection ID to both operators. This makes both operators as dynamic execution engine for queries and not dedicated to load only one table.		
<div><div>✓</div><div>The DAG allows to switch between append-only and delete-load functionality</div></div>		
You did very well on this by using the truncate_table flag, according to which you decide to truncate prior insertion or not based on the loading mode required; append or truncate-insert		
<h3>Data Quality Checks</h3>		
<div><div>✓</div><div>Data quality check is done with correct operator</div></div>		
Correct! you created an operator for data quality to run after loading is finished.		
<div><div>✓</div><div>The DAG either fails or retries n times</div></div>		
Correct! exceptions are raised in case the expected result was not met, this will fail the task and will lead to retry logic as per the config passed to the dag.		
<div><div>✓</div><div>Operator uses params to get the tests and the results, tests are not hard coded to the operator</div></div>		
<p>Excellent! Exactly as required. You passed both tests and expected results as parameters to the operator. This makes your operator dynamic and able to run any tests in a dynamic manner. You can also add another field in the dict you are passing which is the comparison operator; according to which you decide how to compare between the result generated and expected result.</p> <p>For example</p> <pre>dq_checks=[{'check_sql': "SELECT COUNT(*) FROM users WHERE userid is null", 'expected_result': 0, comaprison:'='}, {'check_sql': "SELECT COUNT(*) FROM songs WHERE songid is null", 'expect ed_result': 0 , comparison:'>'}]</pre>		
↓ DOWNLOAD PROJECT		