

Data Engineering 101 - Getting Started with Apache Airflow

[BEGINNER](#)[DATA ENGINEERING](#)[PYTHON](#)

Overview

- Understanding the need for Apache Airflow and its components
- We will create our first DAG to get live cricket scores using Apache Airflow
-

Introduction

Automation of work plays a key role in any industry and it is one of the quickest ways to reach functional efficiency. But many of us fail to understand how to automate some tasks and end in the loop of manually doing the same things again and again.

1998, 1999, 2000, 2001, 2002, 2003, 2004, 2005, 2006, 2007, 2008, 2009, 2010, 2011, 2012, 2013, 2014, 2015, 2016, 2017, 2018, 2019, 2020, 2021, 2022, 2023, 2024, 2025, 2026, 2027, 2028, 2029, 2030, 2031, 2032, 2033, 2034, 2035, 2036, 2037, 2038, 2039, 2040, 2041, 2042, 2043, 2044, 2045, 2046, 2047, 2048, 2049, 2050, 2051, 2052, 2053, 2054, 2055, 2056, 2057, 2058, 2059, 2060, 2061, 2062, 2063, 2064, 2065, 2066, 2067, 2068, 2069, 2070, 2071, 2072, 2073, 2074, 2075, 2076, 2077, 2078, 2079, 2080, 2081, 2082, 2083, 2084, 2085, 2086, 2087, 2088, 2089, 2090, 2091, 2092, 2093, 2094, 2095, 2096, 2097, 2098, 2099, 2100, 2101, 2102, 2103, 2104, 2105, 2106, 2107, 2108, 2109, 2110, 2111, 2112, 2113, 2114, 2115, 2116, 2117, 2118, 2119, 2120, 2121, 2122, 2123, 2124, 2125, 2126, 2127, 2128, 2129, 2130, 2131, 2132, 2133, 2134, 2135, 2136, 2137, 2138, 2139, 2140, 2141, 2142, 2143, 2144, 2145, 2146, 2147, 2148, 2149, 2150, 2151, 2152, 2153, 2154, 2155, 2156, 2157, 2158, 2159, 2160, 2161, 2162, 2163, 2164, 2165, 2166, 2167, 2168, 2169, 2170, 2171, 2172, 2173, 2174, 2175, 2176, 2177, 2178, 2179, 2180, 2181, 2182, 2183, 2184, 2185, 2186, 2187, 2188, 2189, 2190, 2191, 2192, 2193, 2194, 2195, 2196, 2197, 2198, 2199, 2200, 2201, 2202, 2203, 2204, 2205, 2206, 2207, 2208, 2209, 2210, 2211, 2212, 2213, 2214, 2215, 2216, 2217, 2218, 2219, 2220, 2221, 2222, 2223, 2224, 2225, 2226, 2227, 2228, 2229, 2230, 2231, 2232, 2233, 2234, 2235, 2236, 2237, 2238, 2239, 2240, 2241, 2242, 2243, 2244, 2245, 2246, 2247, 2248, 2249, 2250, 2251, 2252, 2253, 2254, 2255, 2256, 2257, 2258, 2259, 2260, 2261, 2262, 2263, 2264, 2265, 2266, 2267, 2268, 2269, 2270, 2271, 2272, 2273, 2274, 2275, 2276, 2277, 2278, 2279, 2280, 2281, 2282, 2283, 2284, 2285, 2286, 2287, 2288, 2289, 2290, 2291, 2292, 2293, 2294, 2295, 2296, 2297, 2298, 2299, 2300, 2301, 2302, 2303, 2304, 2305, 2306, 2307, 2308, 2309, 2310, 2311, 2312, 2313, 2314, 2315, 2316, 2317, 2318, 2319, 2320, 2321, 2322, 2323, 2324, 2325, 2326, 2327, 2328, 2329, 2330, 2331, 2332, 2333, 2334, 2335, 2336, 2337, 2338, 2339, 2340, 2341, 2342, 2343, 2344, 2345, 2346, 2347, 2348, 2349, 2350, 2351, 2352, 2353, 2354, 2355, 2356, 2357, 2358, 2359, 2360, 2361, 2362, 2363, 2364, 2365, 2366, 2367, 2368, 2369, 2370, 2371, 2372, 2373, 2374, 2375, 2376, 2377, 2378, 2379, 2380, 2381, 2382, 2383, 2384, 2385, 2386, 2387, 2388, 2389, 2390, 2391, 2392, 2393, 2394, 2395, 2396, 2397, 2398, 2399, 2400, 2401, 2402, 2403, 2404, 2405, 2406, 2407, 2408, 2409, 2410, 2411, 2412, 2413, 2414, 2415, 2416, 2417, 2418, 2419, 2420, 2421, 2422, 2423, 2424, 2425, 2426, 2427, 2428, 2429, 2430, 2431, 2432, 2433, 2434, 2435, 2436, 2437, 2438, 2439, 2440, 2441, 2442, 2443, 2444, 2445, 2446, 2447, 2448, 2449, 2450, 2451, 2452, 2453, 2454, 2455, 2456, 2457, 2458, 2459, 2460, 2461, 2462, 2463, 2464, 2465, 2466, 2467, 2468, 2469, 2470, 2471, 2472, 2473, 2474, 2475, 2476, 2477, 2478, 2479, 2480, 2481, 2482, 2483, 2484, 2485, 2486, 2487, 2488, 2489, 2490, 2491, 2492, 2493, 2494, 2495, 2496, 2497, 2498, 2499, 2500, 2501, 2502, 2503, 2504, 2505, 2506, 2507, 2508, 2509, 2510, 2511, 2512, 2513, 2514, 2515, 2516, 2517, 2518, 2519, 2520, 2521, 2522, 2523, 2524, 2525, 2526, 2527, 2528, 2529, 2530, 2531, 2532, 2533, 2534, 2535, 2536, 2537, 2538, 2539, 2540, 2541, 2542, 2543, 2544, 2545, 2546, 2547, 2548, 2549, 2550, 2551, 2552, 2553, 2554, 2555, 2556, 2557, 2558, 2559, 2560, 2561, 2562, 2563, 2564, 2565, 2566, 2567, 2568, 2569, 2570, 2571, 2572, 2573, 2574, 2575, 2576, 2577, 2578, 2579, 2580, 2581, 2582, 2583, 2584, 2585, 2586, 2587, 2588, 2589, 2590, 2591, 2592, 2593, 2594, 2595, 2596, 2597, 2598, 2599, 2600, 2601, 2602, 2603, 2604, 2605, 2606, 2607, 2608, 2609, 2610, 2611, 2612, 2613, 2614, 2615, 2616, 2617, 2618, 2619, 2620, 2621, 2622, 2623, 2624, 2625, 2626, 2627, 2628, 2629, 2630, 2631, 2632, 2633, 2634, 2635, 2636, 2637, 2638, 2639, 2640, 2641, 2642, 2643, 2644, 2645, 2646, 2647, 2648, 2649, 2650, 2651, 2652, 2653, 2654, 2655, 2656, 2657, 2658, 2659, 2660, 2661, 2662, 2663, 2664, 2665, 2666, 2667, 2668, 2669, 2670, 2671, 2672, 2673, 2674, 2675, 2676, 2677, 2678, 2679, 26



Most of us have to deal with different workflows like collecting data from multiple databases, preprocessing it, upload it, and report it. Consequently, it would be great if our daily tasks just automatically trigger on defined time, and all the processes get executed in order. Apache Airflow is one such tool that can be very helpful for you. Whether you are Data Scientist, Data Engineer, or Software Engineer you will definitely find this tool useful.

In this article, we will discuss Apache Airflow, how to install it and we will create a sample workflow and code it in Python.

Table of Contents

1. What is Apache Airflow?
2. Features of Apache Airflow
3. Installation Steps
4. Components of Apache Airflow
 1. Webserver
 2. Scheduler
 3. Executor
 4. Metabase
5. User Interface
6. Define your first DAG
7. End Notes

What is Apache Airflow?

[Apache Airflow](#) is a workflow engine that will easily schedule and run your complex data pipelines. It will make sure that each task of your data pipeline will get executed in the correct order and each task gets the required resources.

It will provide you an amazing user interface to monitor and fix any issues that may arise.

Features of Apache Airflow

1. **Easy to Use:** If you have a bit of python knowledge, you are good to go and deploy on Airflow.
2. **Open Source:** It is free and open-source with a lot of active users.
3. **Robust Integrations:** It will give you ready to use operators so that you can work with Google Cloud Platform, Amazon AWS, Microsoft Azure, etc.
4. **Use Standard Python to code:** You can use python to create simple to complex workflows with complete flexibility.
5. **Amazing User Interface:** You can monitor and manage your workflows. It will allow you to check the status of completed and ongoing tasks.

Installation Steps

Let's start with the installation of the Apache Airflow. Now, if already have pip installed in your system, you can skip the first command. To install pip run the following command in the terminal.

```
sudo apt-get install python3-pip
```

Next airflow needs a home on your local system. By default **~/airflow** is the default location but you can change it as per your requirement.

```
export AIRFLOW_HOME=~/airflow
```

Now, install the apache airflow using the pip with the following command.

```
pip3 install apache-airflow
```

Airflow requires a database backend to run your workflows and to maintain them. Now, to initialize the database run the following command.

```
airflow initdb
```

We have already discussed that airflow has an amazing user interface. To start the webserver run the following command in the terminal. The default port is 8080 and if you are using that port for something else then you can change it.

```
airflow webserver -p 8080
```

Now, start the airflow scheduler using the following command in a different terminal. It will run all the time and monitor all your workflows and triggers them as you have assigned.

```
airflow scheduler
```

Now, create a folder name dags in the airflow directory where you will define your workflows or DAGs and open the web browser and go open: <http://localhost:8080/admin/> and you will see something like this:

Components of Apache Airflow

- **DAG:** It is the Directed Acyclic Graph – a collection of all the tasks that you want to run which is organized and shows the relationship between different tasks. It is defined in a python script.
- **Web Server:** It is the user interface built on the Flask. It allows us to monitor the status of the DAGs and trigger them.
- **Metadata Database:** Airflow stores the status of all the tasks in a database and do all read/write operations of a workflow from here.
- **Scheduler:** As the name suggests, this component is responsible for scheduling the execution of DAGs. It retrieves and updates the status of the task in the database.

User Interface

Now that you have installed the Airflow, let's have a quick overview of some of the components of the user interface.

DAGS VIEW

It is the default view of the user interface. This will list down all the DAGS present in your system. It will give you a summarized view of the DAGS like how many times a particular DAG was run successfully, how many times it failed, the last execution time, and some other useful links.

GRAPH VIEW

In the graph view, you can visualize each and every step of your workflow with their dependencies and their current status. You can check the current status with different color codes like:

TREE VIEW

The tree view also represents the DAG. If you think your pipeline took a longer time to execute than expected then you can check which part is taking a long time to execute and then you can work on it.

TASK DURATION

In this view, you can compare the duration of your tasks run at different time intervals. You can optimize your algorithms and compare your performance here.

CODE

In this view, you can quickly view the code that was used to generate the DAG.

Define your first DAG

Let's start and define our first DAG.

In this section, we will create a workflow in which the first step will be to print “Getting Live Cricket Scores” on the terminal, and then using an API, we will print the live scores on the terminal. Let’s test the API first and for that, you need to install the **cricket-cli** library using the following command.

```
sudo pip3 install cricket-cli
```

Now, run the following command and get the scores.

```
cricket scores
```

It might take a few seconds of time, based on your internet connection, and will return you the output something like this:

Importing the Libraries

Now, we will create the same workflow using Apache Airflow. The code will be completely in python to define a DAG. Let’s start with importing the libraries that we need. We will use only the BashOperator only as our workflow requires the Bash operations to run only.

```
1 from datetime import timedelta
2
3 # The DAG object; we'll need this to instantiate a DAG
4 from airflow import DAG
5 # Operators; we need this to operate!
6 from airflow.operators.bash_operator import BashOperator
7 from airflow.utils.dates import days_ago
```

airflow_import.py hosted with ❤ by GitHub

[view raw](#)

Defining DAG Arguments

For each of the DAG, we need to pass one argument dictionary. Here is the description of some of the arguments that you can pass:

- **owner**: The name of the owner of the workflow, should be alphanumeric and can have underscores but should not contain any spaces.

- **depends_on_past:** If each time you run your workflow, the data depends upon the past run then mark it as True otherwise mark it as False.
- **start_date:** Start date of your workflow
- **email:** Your email ID, so that you can receive an email whenever any task fails due to any reason.
- **retry_delay:** If any task fails, then how much time it should wait to retry it.

```

1  # You can override them on a per-task basis during operator initialization
2  default_args = {
3      'owner': 'lakshay',
4      'depends_on_past': False,
5      'start_date': days_ago(2),
6      'email': ['airflow@example.com'],
7      'email_on_failure': False,
8      'email_on_retry': False,
9      'retries': 1,
10     'retry_delay': timedelta(minutes=5),
11     # 'queue': 'bash_queue',
12     # 'pool': 'backfill',
13     # 'priority_weight': 10,
14     # 'end_date': datetime(2016, 1, 1),
15     # 'wait_for_downstream': False,
16     # 'dag': dag,
17     # 'sla': timedelta(hours=2),
18     # 'execution_timeout': timedelta(seconds=300),
19     # 'on_failure_callback': some_function,
20     # 'on_success_callback': some_other_function,
21     # 'on_retry_callback': another_function,
22     # 'sla_miss_callback': yet_another_function,
23     # 'trigger_rule': 'all_success'
24 }

```

[view raw](#)

args_airflow.py hosted with ❤ by GitHub

Defining DAG

Now, we will create a DAG object and pass the **dag_id** which is the name of the DAG and it should be unique. Pass the arguments that we defined in the last step and add a description and **schedule_interval** which will run the DAG after the specified interval of time

```

1  # define the DAG
2  dag = DAG(
3      'live_cricket_scores',
4      default_args=default_args,
5      description='First example to get Live Cricket Scores',
6      schedule_interval=timedelta(days=1),
7  )

```

[view raw](#)

dag_airflow.py hosted with ❤ by GitHub

Defining the Tasks

We will have 2 tasks for our workflow:

1. **print:** In the first task, we will print the “Getting Live Cricket Scores!!!” on the terminal using the echo command.
2. **get_cricket_scores:** In the second task, we will print the live cricket scores using the library that we have installed.

Now, while defining the task first we need to choose the right operator for the task. Here both the commands are terminal-based so we will use the BashOperator.

We will pass the **task_id** which is a unique identifier of the task and you will see this name on the nodes of Graph View of your DAG. Pass the bash command that you want to run and finally the DAG object to which you want to link this task.

Finally, create the pipeline by adding the “>>” operator between the tasks.

```
1  # define the first task
2  t1 = BashOperator(
3      task_id='print',
4      bash_command='echo Getting Live Cricket Scores!!!',
5      dag=dag,
6  )
7
8
9  # define the second task
10 t2 = BashOperator(
11     task_id='get_cricket_scores',
12     bash_command='cricket scores',
13     dag=dag,
14 )
15
16 # task pipeline
17 t1 >> t2
```

define_task.py hosted with ❤ by GitHub [view raw](#)

Update the DAGS in Web UI

Now, refresh the user interface and you will see your DAG in the list. Turn on the toggle on the left of each of the DAG and then trigger the DAG.

Click on the DAG and open the graph view and you will see something like this. Each of the steps in the workflow will be in a separate box and its border will turn dark green once it is completed successfully.

Click on the node “get_cricket scores” to get more details about this step. You will see something like this.

Now, click on View Log to see the output of your code.

That's it. You have successfully created your first DAG in the Apache Airflow.

End Notes

In this article, we have seen the features of Apache Airflow, its user interface components and we have created a simple DAG. In the upcoming article, we will discuss some more concepts like variables, branching, and will create a more complex workflow.

I recommend you go through the following data engineering resources to enhance your knowledge-

- [Getting Started with Apache Hive – A Must Know Tool For all Big Data and Data Engineering Professionals](#)
- [Introduction to the Hadoop Ecosystem for Big Data and Data Engineering](#)
- [Types of Tables in Apache Hive – A Quick Overview](#)

If you have any questions related to this article do let me know in the comments section below.

Article Url - <https://www.analyticsvidhya.com/blog/2020/11/getting-started-with-apache-airflow/>



LAKSHAY ARORA

Ideas have always excited me. The fact that we could dream of something and bring it to reality fascinates me. Computer Science provides me a window to do exactly that. I love programming and use it to solve problems and a beginner in the field of Data Science.