

# Airflow: Lesser Known Tips, Tricks, and Best Practises



Kaxil Naik

Follow

Dec 25, 2018 · 5 min read

There are certain things with all the tools you use that you won't know even after using it for a long time. And once you know it you are like "I wish I knew this before" as you had already told your client that it can't be done in any better way 🤔🤔. **Airflow** like other tool is no different, there are some hidden gems that can make your life easy and make **DAG** development fun.

You might already know some of them and if you know them all — well you are a PRO then 🎩👑.



**Airflow: Lesser Known Tips, Tricks & Best Practises**

# (1) DAG with context Manager

Were you annoyed with yourself when you forgot to add `dag=dag` to your task and Airflow error'ed? Yes, it is easy to forget adding it for each task. It is also redundant to add the same parameter as shown in the following example ( `example_dag.py` file):

```
1  # Normal DAG without Context Manager
2  args = {
3      'owner': 'airflow',
4      'start_date': airflow.utils.dates.days_ago(2),
5  }
6
7  dag = DAG(
8      dag_id='example_dag',
9      default_args=args,
10     schedule_interval='0 0 * * *',
11 )
12
13 run_this_last = DummyOperator(
14     task_id='run_this_last',
15     dag=dag, # You need to repeat this for each task
16 )
17
18 run_this_first = BashOperator(
19     task_id='run_this_first',
20     bash_command='echo 1',
21     dag=dag, # You need to repeat this for each task
22 )
23
24 run_this_first >> run_this_last
```

`example_dag.py` hosted with ❤️ by GitHub

[view raw](#)

```
1  # DAG with Context Manager
2
3  args = {
4      'owner': 'airflow',
5      'start_date': airflow.utils.dates.days_ago(2),
6  }
7
8  with DAG(dag_id='example_dag', default_args=args, schedule_interval='0 0 * * *') as dag:
9
10     run_this_last = DummyOperator(
11         task_id='run_this_last'
12     )
```

```
--
13
14     run_this_first = BashOperator(
15         task_id='run_this_first',
16         bash_command='echo 1'
17     )
18
19     run_this_first >> run_this_last
```

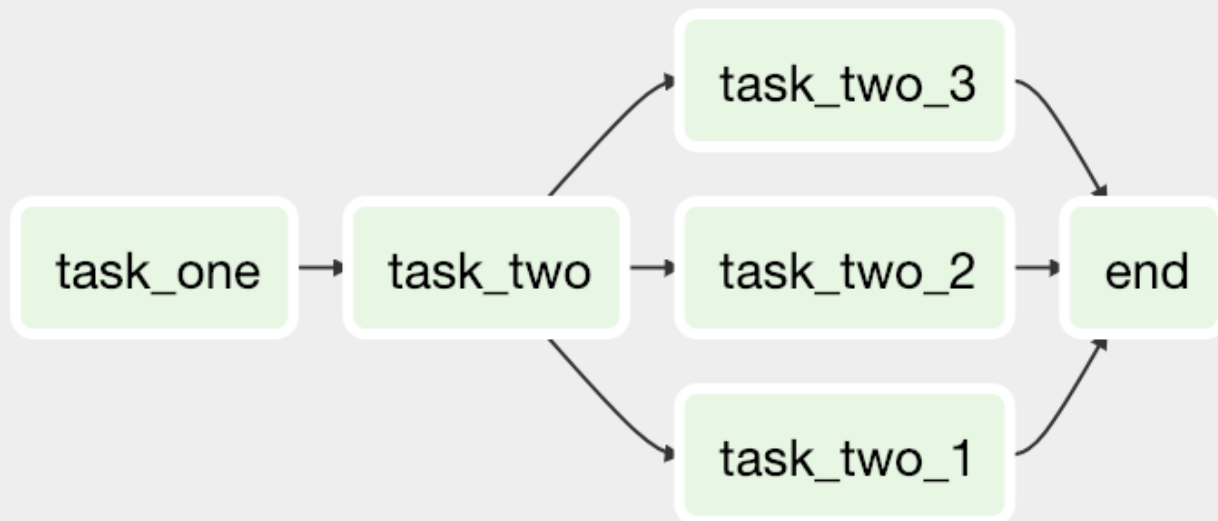
example\_dag\_with\_context.py hosted with ❤ by GitHub

[view raw](#)

The example ( `example_dag.py` file) above just has 2 tasks, but if you have 10 or more then the redundancy becomes more evident. To avoid this you can use **Airflow** DAGs as **context managers** to automatically assign new operators to that DAG as shown in the above example ( `example_dag_with_context.py` ) using `with` statement.

## (2) Using List to set Task dependencies

When you want to create the DAG similar to the one shown in the image below, you would have to repeat task names when setting task dependencies.



```
1 # Setting task dependencies (the NORMAL way)
2 task_one >> task_two
```

```

3 task_two >> task_two_1 >> end
4 task_two >> task_two_2 >> end
5 task_two >> task_two_3 >> end
6
7 # Using Lists (being a PRO :-D )
8 task_one >> task_two >> [task_two_1, task_two_2, task_two_3] >> end

```

airflow\_list\_task\_dependencies.py hosted with ❤ by GitHub

[view raw](#)

As shown in the above code snippet, using our normal way of setting task dependencies would mean that `task_two` and `end` are repeated 3 times. This can be replaced using python lists to achieve the same result in a more elegant way.

### (3) Use default arguments to avoid repeating arguments

Airflow allowing passing a dictionary of parameters that would be available to all the task in that DAG.

For example, at DataReply, we use BigQuery for all our DataWarehouse related DAGs and instead of passing parameters like `labels`, `bigquery_conn_id` to each task, we simply pass it in `default_args` dictionary as shown in the DAG below.

```

1  default_args = {
2      'owner': 'airflow',
3      'depends_on_past': False,
4      'start_date': airflow.utils.dates.days_ago(2),
5      # All the parameters below are BigQuery specific and will be available to all the ta
6      'bigquery_conn_id': 'gcp-bigquery-connection',
7      'write_disposition': 'WRITE_EMPTY',
8      'create_disposition': 'CREATE_IF_NEEDED',
9      'labels': {'client': 'client-1'}
10 }
11
12 with DAG(dag_id='airflow_tutorial_gcp', default_args=default_args, schedule_interval=Nor
13
14     query_1 = BigQueryOperator(
15         task_id='query_1',
16         sql='select 1'
17     )
18
19     query_2 = BigQueryOperator(
20         task_id='query_2',

```

```
21         sql='select 1'
22     )
23
24     query_1 >> query_2
```

airflow\_default\_args.py hosted with ❤ by GitHub

[view raw](#)

This is also useful when you want alerts on individual task failures instead of just DAG failures which I already mentioned in my last blog post on Integrating Slack Alerts in Airflow.

## (4) The "params" argument

“params” is a dictionary of DAG level parameters that are made accessible in templates. These params can be overridden at the task level.

This is an extremely helpful argument and I have been personally using it a lot as it can be accessed in **templated** field with jinja templating using `params.param_name`. An example usage is as follows:

```
1  # You can pass `params` dict to DAG object
2  default_args = {
3      'owner': 'airflow',
4      'depends_on_past': False,
5      'start_date': airflow.utils.dates.days_ago(2),
6  }
7
8  dag = DAG(
9      dag_id='airflow_tutorial_2',
10     default_args=default_args,
11     schedule_interval=None,
12     params={
13         "param1": "value1",
14         "param2": "value2"
15     }
16 )
17
18 bash = BashOperator(
19     task_id='bash',
20     bash_command='echo {{ params.param1 }}', # Output: value1
21     dag=dag
22 )
```

&lt;&lt; )

airflow\_params\_usage\_1.py hosted with ❤ by GitHub

[view raw](#)

```
1  # Passing `params` dict in `default_arg` dict
2  default_args = {
3      'owner': 'airflow',
4      'depends_on_past': False,
5      'start_date': airflow.utils.dates.days_ago(2),
6      'params': {
7          "param1": "value2",
8          "param2": "value1"
9      }
10 }
11
12 dag = DAG(
13     dag_id='airflow_tutorial_2',
14     default_args=default_args,
15     schedule_interval=None,
16 )
17
18 bash = BashOperator(
19     task_id='bash',
20     bash_command='echo {{ params.param1 }}', # Output: value2
21     dag=dag
22 )
```

airflow\_params\_usage\_2.py hosted with ❤ by GitHub

[view raw](#)

```
1  # Passing `params` dict in tasks
2
3  default_args = {
4      'owner': 'airflow',
5      'depends_on_past': False,
6      'start_date': airflow.utils.dates.days_ago(2),
7  }
8
9  dag = DAG(
10     dag_id='airflow_tutorial_2',
11     default_args=default_args,
12     schedule_interval=None,
13 )
14
15 bash = BashOperator(
16     task_id='bash',
17     bash_command='echo {{ params.param1 }}', # Output: value3
18     params={
```

```

19         "param1": "value3",
20         "param2": "value4"
21     }
22     dag=dag
23 )

```

airflow\_params\_usage\_3.py hosted with ❤ by GitHub

[view raw](#)

```

1  # You can override `params` dict passed in DAG object in `default_arg` dict
2  default_args = {
3      'owner': 'airflow',
4      'depends_on_past': False,
5      'start_date': airflow.utils.dates.days_ago(2),
6      'params': {
7          "param1": "value2",
8          "param2": "value2"
9      }
10 }
11
12 dag = DAG(
13     dag_id='airflow_tutorial_2',
14     default_args=default_args,
15     schedule_interval=None,
16     params={
17         "param1": "value1",
18         "param2": "value2"
19     }
20 )
21
22 # You can override `params` dict passed in DAG object or `default_arg` in each individual task
23 bash = BashOperator(
24     task_id='bash',
25     bash_command='echo {{ params.param1 }}', # Output: value3
26     params={
27         "param1": "value3"
28     }
29     dag=dag
30 )

```

airflow\_params\_usage\_4.py hosted with ❤ by GitHub

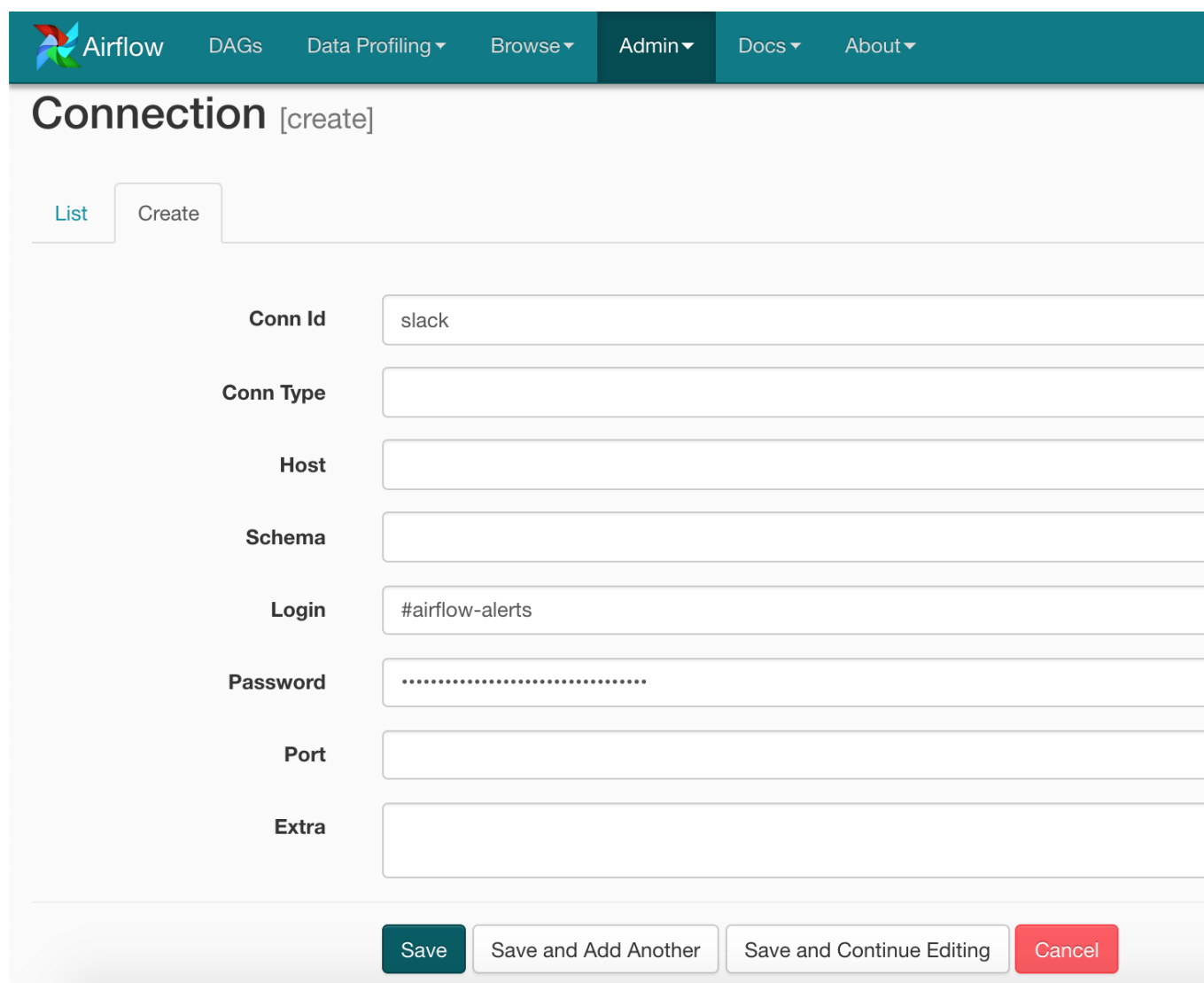
[view raw](#)

It makes it easy for you to write parameterized DAG instead of hard-coding values. Also as shown in the examples above `params` dictionary can be defined at 3 places: (1) In DAG object (2) In `default_args` dictionary (3) Each task.

## (5) Storing Sensitive data in Connections

Most users are aware of this but I have still seen passwords stored in plain-text inside the DAG. For goodness sake — don't do that. You should write your DAGs in a way that you are confident enough to store your DAGs in a public repository.

By default, Airflow will save the passwords for the connection in plain text within the metadata database. The `crypto` package is highly recommended during Airflow installation and can be simply done by `pip install apache-airflow[crypto]`.



The screenshot shows the Airflow web interface with a teal header bar containing the Airflow logo and navigation links: DAGs, Data Profiling, Browse, Admin, Docs, and About. Below the header, the main content area is titled 'Connection [create]'. There are two tabs: 'List' (inactive) and 'Create' (active). The 'Create' tab contains a form with the following fields: 'Conn Id' (text input with 'slack'), 'Conn Type' (text input), 'Host' (text input), 'Schema' (text input), 'Login' (text input with '#airflow-alerts'), 'Password' (password input with masked dots), 'Port' (text input), and 'Extra' (text input). At the bottom of the form are four buttons: 'Save' (teal), 'Save and Add Another' (light gray), 'Save and Continue Editing' (light gray), and 'Cancel' (red).

You can then easily access it as follows:

```
from airflow.hooks.base_hook import BaseHook
slack_token = BaseHook.get_connection('slack').password
```



## (6) Restrict the number of Airflow variables in your DAG







Airflow Variables are stored in Metadata Database, so any call to variables would mean a connection to Metadata DB. Your DAG files are parsed every X seconds. Using a large number of variable in your DAG (and worse in `default_args` ) may mean you might end up saturating the number of allowed connections to your database.

To avoid this situation, you can either just use a single Airflow variable with JSON value. As an Airflow variable can contain JSON value, you can store all your DAG configuration inside a single variable as shown in the image below:

### Variables

No file chosen

List (3)

<input type="checkbox"/>		Key	Val
<input type="checkbox"/>	 	dag1_config	{"var1 ":"value1","var2 ":"value2"}
<input type="checkbox"/>	 	var1	value1
<input type="checkbox"/>	 	var2	value2

As shown in this screenshot you can either store values in separate Airflow variables or under a single Airflow variable as a JSON field

You can then access them as shown below under **Recommended** way:

```
1 from airflow.models import Variable
2
3 # Common (Not-so-nice way)
4 # 3 DB connections when the file is parsed
5 var1 = Variable.get("var1")
6 var2 = Variable.get("var2")
7 var3 = Variable.get("var3")
8
```

```

9  # Recommended Way
10 # Just 1 Database call
11 dag_config = Variable.get("dag1_config", deserialize_json=True)
12 dag_config["var1"]
13 dag_config["var2"]
14 dag_config["var3"]
15
16 # You can directly use it Templated arguments {{ var.json.my_var.path }}
17 bash_task = BashOperator(
18     task_id="bash_task",
19     bash_command='{{ var.json.dag1_config.var1 }} ',
20     dag=dag,
21 )

```

airflow\_json\_variables.py hosted with ❤ by GitHub

[view raw](#)

## (7) The "context" dictionary

Users often forget the contents of the `context` dictionary when using `PythonOperator` with a callable function.

The context contains references to related objects to the task instance and is documented under the macros section of the API as they are also available to templated field.

```

{
    'dag': task.dag,
    'ds': ds,
    'next_ds': next_ds,
    'next_ds_nodash': next_ds_nodash,
    'prev_ds': prev_ds,
    'prev_ds_nodash': prev_ds_nodash,
    'ds_nodash': ds_nodash,
    'ts': ts,
    'ts_nodash': ts_nodash,
    'ts_nodash_with_tz': ts_nodash_with_tz,
    'yesterday_ds': yesterday_ds,
    'yesterday_ds_nodash': yesterday_ds_nodash,
    'tomorrow_ds': tomorrow_ds,
    'tomorrow_ds_nodash': tomorrow_ds_nodash,
    'END_DATE': ds,
    'end_date': ds,
    'dag_run': dag_run,
    'run_id': run_id,
    'execution_date': self.execution_date,
    'prev_execution_date': prev_execution_date,

```

```

'next_execution_date': next_execution_date,
'latest_date': ds,
'macros': macros,
'params': params,
'tables': tables,
'task': task,
'task_instance': self,
'ti': self,
'task_instance_key_str': ti_key_str,
'conf': configuration,
'test_mode': self.test_mode,
'var': {
    'value': VariableAccessor(),
    'json': VariableJsonAccessor()
},
'inlets': task.inlets,
'outlets': task.outlets,
}

```

## (8) Generating Dynamic Airflow Tasks

I have been answering many questions on StackOverflow on how to create dynamic tasks. The answer is simple, you just need to generate **unique** `task_id` for all of your tasks. Below are 2 examples on how to achieve that:

```

1  # Using DummyOperator
2  a = []
3  for i in range(0,10):
4      a.append(DummyOperator(
5          task_id='Component'+str(i),
6          dag=dag))
7      if i != 0:
8          a[i-1] >> a[i]
9
10 # From a List
11 sample_list = ["val1", "val2", "val3"]
12 tasks_list = []
13 for index, value in enumerate(sample_list):
14     tasks_list.append(DummyOperator(
15         task_id='Component'+str(index),
16         dag=dag))
17     if index != 0:
18         tasks_list[index-1] >> tasks_list[index]

```

airflow\_dynamic\_task.py hosted with ❤️ by GitHub

[view raw](#)

## (9) Run "airflow upgradedb" instead of "airflow initdb"

Thanks to Ash Berlin for this tip in his talk in the First Apache Airflow London Meetup.

`airflow initdb` will create all default connections, charts etc that we might not use and don't want in our production database. `airflow upgradedb` will instead just apply any missing migrations to the database table. (including creating missing tables etc.) It is also safe to run every time, it tracks which migrations have already been applied (using the Alembic module).

. . .

Let me know in the comments section below if you know something that would be worth adding in this blog post. Happy Airflow'ing :-)

[Python](#) [Apache Airflow](#) [Airflow](#)

[About](#) [Help](#) [Legal](#)