# Airflow详细搭建过程（亲测 + 总结）-CSDN博客

　Airflow 专栏收录该内容

2 篇文章 0 订阅

订阅专栏

## Airflow详细搭建过程（亲测+总结）

官网:Apache airflow
Airflow是社区创建的一个平台，用于以编程方式编写，安排和监视工作流。
半个月前搭了一次，现在做一个记录，不逼逼，开始搭建了，全程多图：

### 环境准备

系统: cent os 7
conda版本: 4.8.2
airflow版本 1.10.11

### 开始搭建

我这边会用一个conda创建一个apache airflow的环境:

```
conda create -n airflow_env python=3.7
    1
```

切换到当前的这个环境：

```
conda activate airflow_env
    1
```

### 搭建airflow

搭建airflow的话，官网有一套详细的文档Airflow 搭建
接下来就按照这个方式来

```
# airflow needs a home, ~/airflow is the default,
# but you can lay foundation somewhere else if you prefer
# (optional)
export AIRFLOW_HOME=~/airflow

# install from pypi using pip
pip install apache-airflow

# initialize the database
airflow initdb

# start the web server, default port is 8080
# airflow webserver -p 8080     这里做个修改，后面加上-D参数让它后台运行
airflow webserver -p 8080 -D

# start the scheduler
airflow scheduler

# visit localhost:8080 in the browser and enable the example dag in the home page
```
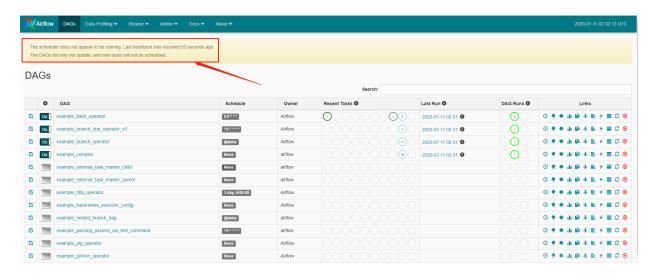
经过上述的步骤的话，可以在浏览器中打开这个页面:



搭建到这里，可以看到一个基础的样子，但是从目前而言，我能发现的有以下几个问题：

- airflow的元数据存储信息默认是使用`sqlite`进行存储。
- 页面的右上角的时间是比正常的时间晚了8个小时的

- airflow默认是用单线程调度任务的，如下图（之后就知道了）



现在问题已经暴露了，那么就开始处理吧。

airflow的元数据存储默认用`sqlite`，现在切换为`mysql`

`sqlite`不支持多线程，所以我打算切换成`mysql`,`mysql`的安装可以参考:Linux centos安装mysql

`airflow`中可选很多其它的选项，那么在`airflow`里面我们可以选择一些所需的组件，例如现在需要利用`mysql`存储`airflow`的相关信息，可选的功能列表在此处airlfow的其它选项

就需要`pip install 'apache-airflow[mysql]'`

这里，我就不挑了，直接来全套插件,这样就会遇到更多的坑：

```
yum install mysql-devel gcc gcc-devel python-devel krb5-devel.x86_64 cyrus-sasl-devel -y
pip install 'apache-airflow[all]' -i https://pypi.tuna.tsinghua.edu.cn/simple/
```

- 1
- 2

这个时候可能会遇到很多的错误，例如下面的:



这里面的错误要有点耐心去处理一下，在这里面一般都是版本上的错误，将红字里面出现的版本卸载掉，重新指定版本号安装就好了。

## 配置airflow

- 配置`mysql`作为`airflow`的元数据信息的存储：
  - `vim /etc/my.cnf` 在`[mysqld]`下面添加`explicit_defaults_for_timestamp=1`
  - `systemctl restart mysqld`重启`mysql`
  - 在`mysql`中创建一个数据库，我这边的数据库的名字叫做`airflow`

```
mysql> show databases;
+--------------------+
| Database           |
+--------------------+
| information_schema |
| airflow            |
| mysql              |
| performance_schema |
| sys                |
+--------------------+
5 rows in set (0.00 sec)

mysql> █  https://blog.csdn.net/OldDirverHelpMe
```

  - 进入`airflow`的文件夹下，修改`airflow.cfg`。

```
# The SqlAlchemy connection string to the metadata database.
# SqlAlchemy supports many different database engine, more information
# their website
# sql_alchemy_conn = sqlite:////root/airflow/airflow.db
sql_alchemy_conn = mysql://192.168.100.120/airflow?user=root&password=12345678
```

- 解决`airflow`一次只能执行一个任务的问题:

```
# The executor class that airflow should use. Choices include
# SequentialExecutor, LocalExecutor, CeleryExecutor, DaskExecutor, KubernetesExecutor
# executor = SequentialExecutor
executor = LocalExecutor
```

相关的依据请参考<u>airflow 执行者</u>

- 解决`airflow`时间晚了八个小时的问题
  修改`airflow.cfg`的内容:

```
# Default timezone in case supplied date times are naive
# can be utc (default), system, or any IANA timezone string (e.g. Europe/Amsterdam)
# default_timezone = utc
default_timezone = Asia/Shanghai
```

这部分的内容就要参考<u>airflow 修改源码</u>
具体的内容如下:

#找到airflow的安装位置
find / -name airflow

- 1
- 2

```
(base) [root@spark2 ~]# find / -name airflow
/root/anaconda3/pkgs/airflow-1.10.10-py37_0/lib/python3.7/site-packages/airflow
/root/anaconda3/pkgs/airflow-1.10.10-py37_0/lib/python3.7/site-packages/airflow/www_rbac/templates/airflow
/root/anaconda3/pkgs/airflow-1.10.10-py37_0/lib/python3.7/site-packages/airflow/www/templates/airflow
/root/anaconda3/pkgs/airflow-1.10.10-py37_0/lib/python3.7/site-packages/airflow/bin/airflow
/root/anaconda3/pkgs/airflow-1.10.10-py37_0/bin/airflow
/root/anaconda3/envs/ame/lib/python3.7/site-packages/airflow
/root/anaconda3/envs/ame/lib/python3.7/site-packages/airflow/bin/airflow
/root/anaconda3/envs/ame/lib/python3.7/site-packages/airflow/www/templates/airflow
/root/anaconda3/envs/ame/lib/python3.7/site-packages/airflow/www_rbac/templates/airflow
/root/anaconda3/envs/ame/bin/airflow
/root/anaconda3/envs/airflow_env/bin/airflow
/root/anaconda3/envs/airflow_env/lib/python3.7/site-packages/airflow
/root/anaconda3/envs/airflow_env/lib/python3.7/site-packages/airflow/bin/airflow
/root/anaconda3/envs/airflow_env/lib/python3.7/site-packages/airflow/www/templates/airflow
/root/anaconda3/envs/airflow_env/lib/python3.7/site-packages/airflow/www_rbac/templates/airflow
/root/airflow
/root/airflow/logs/anaconda3/envs/airflow_env/lib/python3.7/site-packages/airflow
```
https://blog.csdn.net/OldDirverHelpMe

1. 修改`/root/anaconda3/envs/airflow_env/lib/python3.7/site-packages/airflow/utils/timezone.py`

   在第27行下面添加:

```
26 # UTC time zone as a tzinfo instance.
27 utc = pendulum.timezone('UTC')
28 from airflow import configuration as conf
29 try:
30         tz = conf.get("core", "default_timezone")
31         if tz == "system":
32                 utc = pendulum.local_timezone()
33         else:
34                 utc = pendulum.timezone(tz)
35 except Exception:
36     pass
37
38
```

```
from airflow import configuration as conf
try:
        tz = conf.get("core", "default_timezone")
        if tz == "system":
                utc = pendulum.local_timezone()
        else:
                utc = pendulum.timezone(tz)
except Exception:
        pass
```

修改 utcnow() 函数:

```
60
61 def utcnow():
62     """
63     Get the current date and time in UTC
64
65     :return:
66     """
67
68     # pendulum utcnow() is not used as that sets a TimezoneInfo object
69     # instead of a Timezone. This is not pickable and also creates issues
70     # when using replace()
71     #d = dt.datetime.utcnow()
72     d = dt.datetime.now()
73     d = d.replace(tzinfo=utc)
74
75     return d
76
77
```

2. 修改`/root/anaconda3/envs/airflow_env/lib/python3.7/site-packages/airflow/utils/sqlalchemy.py`

   在第38行下面添加: `utc = pendulum.timezone('UTC')`

```
from airflow import configuration as conf
try:
        tz = conf.get("core", "default_timezone")
        if tz == "system":
                utc = pendulum.local_timezone()
        else:
                utc = pendulum.timezone(tz)
except Exception:
        pass
```



注释掉这个:



3. 修改 /root/anaconda3/envs/airflow_env/lib/python3.7/site-packages/airflow/www/templates/admin/master.html
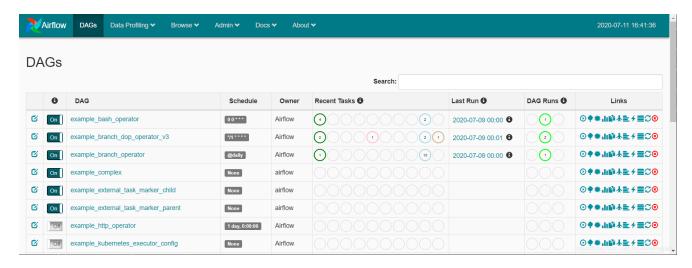
修改下箭头所指的行，修改的内容在注释的下面

此时修改完毕了，这个时候重启一下 `airflow`就好了。修改源码的操作参考:airflow 修改中国时区(改airflow源码)

## 最后的结果如下：

使用了`mysql`作为元数据信息的存储，在`mysql`里面可以看到:

```
Type 'help;' or '\h' for help. Type '\c' to clear the current input statement.

mysql> use airflow
Reading table information for completion of table and column names
You can turn off this feature to get a quicker startup with -A

Database changed
mysql> show tables;
+-----------------------------+
| Tables_in_airflow           |
+-----------------------------+
| alembic_version             |
| chart                       |
| connection                  |
| dag                         |
| dag_code                    |
| dag_pickle                  |
| dag_run                     |
| dag_tag                     |
| import_error                |
| job                         |
| known_event                 |
| known_event_type            |
| kube_resource_version       |
| kube_worker_uuid            |
| log                         |
| rendered_task_instance_fields |
| serialized_dag              |
| sla_miss                    |
| slot_pool                   |
| task_fail                   |
| task_instance               |
| task_reschedule             |
| users                       |
| variable                    |
| xcom                        |
+-----------------------------+
25 rows in set (0.00 sec)
```

登录`airflow`的时候:

解决了之前提到的三个问题