

ETL. Семинар №5

1. Викторина
2. Блок 1

Задание 1

Используя материалы лекции установите Apache Airflow в чистый virtualenv с использованием constraint файла. Убедитесь что по адресу localhost:8080 открывается список Dag сценариев



40 минут

Решение:

Итак, давайте установим Airflow. Вначале создадим новый virtualenv, чтобы не устанавливать в корень системы все библиотеки. После этого выполним следующие команды

```
1 AIRFLOW_VERSION=2.5.0
2 PYTHON_VERSION="$(python --version | cut -d " " -f 2 | cut -d "." -f 1-2)"
3 # For example: 3.7
4 CONSTRAINT_URL="https://raw.githubusercontent.com/apache/airflow/constraints-${AIRFLOW_VERSION}/constraints-no-providers-${PYTHON_VERSION}.txt"
5 # For example: https://raw.githubusercontent.com/apache/airflow/constraints-2.5.0/constraints-no-providers-3.7.txt
6 pip install "apache-airflow==${AIRFLOW_VERSION}" --constraint "${CONSTRAINT_URL}"
```

Мы указываем версию Airflow которую хотим установить, версию python которую будем использовать, ссылку на файл ограничений и выполняем установку с помощью PyPI используя параметры, определенные нами заранее.

Если после установки команда airflow не распознается системой (чаще всего это происходит на Windows при использовании WSL) необходимо убедиться, что в

переменные окружения добавлен путь к установочной директории Airflow и если нет, то добавить.

```
1 PATH=$PATH:~/local/bin
```

Или можно вызывать Airflow с помощью python.

```
1 python -m airflow
```

Также очень часто встречается ошибка при вызове Airflow Symbol not found: _Py_GetArgcArgv. Эта ошибка может означать что вы используете несовместимую версию python. Суть проблемы в том, что библиотека Airflow зависит от setproctitle, использует непубличный API Python, который недоступен в стандартной установке /usr/local/opt/(которая указывает на путь под /usr/local/Cellar).

Простое решение — просто убедиться, что вы используете версию Python, в которой есть dylib библиотеки Python. Например:

```
1 # Note: these instructions are for python3.7 but can be loosely modified for other
  versions
2 brew install python@3.7
3 virtualenv -p
  /usr/local/opt/python@3.7/Frameworks/Python.framework/Versions/3.7/bin/python3 .toy-venv
4 source .toy-venv/bin/activate
5 pip install apache-airflow
6 python
7 >>> import setproctitle
8 # Success!
```

Кроме того, вы можете загрузить и установить Python прямо с официального сайта .

Задание 2

Создайте нового пользователя с правами админа и авторизуйтесь с его помощью в интерфейсе airflow



15 минут

Решение:

Остановить выполнение уже запущенной копии airflow,

```
# create an admin user
airflow users create \
  --username admin \
  --firstname Peter \
  --lastname Parker \
  --role Admin \
  --email spiderman@superhero.org
```

Задание 3

Посмотрите список существующих пайплайнов. Запустите некоторые из них, посмотрите на результат выполнения. Откройте логи выполнения пайплайнов. Отправьте в чат скриншот из логов, говорящий об успешном завершении пайплайна



15 минут

Airflow

[DAGs](#)
[Datasets](#)
[Security](#)
[Browse](#)
[Admin](#)
[Docs](#)

21:33 UTC

DAGs

All 46

Active 0

Paused 0

Auto-refresh

DAG ↕	Owner ↕	Runs ↕	Schedule	Last Run ↕	Next Run ↕ ⌚	Recent Tasks ↕	Actions	Links
<div>dataset_consumes_1</div> <div>consumes dataset-scheduled</div>	airflow	<div> <div></div> <div></div> <div></div> <div></div> </div>	Dataset ⌚		On sb://dag1/output_1.txt	<div> <div></div> <div></div> <div></div> <div></div> <div></div> <div></div> <div></div> <div></div> <div></div> <div></div> </div>	<div>▶</div> <div>⏻</div>	...
<div>dataset_consumes_1_and_2</div> <div>consumes dataset-scheduled</div>	airflow	<div> <div></div> <div></div> <div></div> <div></div> </div>	Dataset ⌚		0 of 2 datasets updated	<div> <div></div> <div></div> <div></div> <div></div> <div></div> <div></div> <div></div> <div></div> <div></div> <div></div> </div>	<div>▶</div> <div>⏻</div>	...
<div>dataset_consumes_1_never_scheduled</div> <div>consumes dataset-scheduled</div>	airflow	<div> <div></div> <div></div> <div></div> <div></div> </div>	Dataset ⌚		0 of 2 datasets updated	<div> <div></div> <div></div> <div></div> <div></div> <div></div> <div></div> <div></div> <div></div> <div></div> <div></div> </div>	<div>▶</div> <div>⏻</div>	...
<div>dataset_consumes_unknown_never_scheduled</div> <div>dataset-scheduled</div>	airflow	<div> <div></div> <div></div> <div></div> <div></div> </div>	Dataset ⌚		0 of 2 datasets updated	<div> <div></div> <div></div> <div></div> <div></div> <div></div> <div></div> <div></div> <div></div> <div></div> <div></div> </div>	<div>▶</div> <div>⏻</div>	...
<div>dataset_produces_1</div> <div>dataset-scheduled produces</div>	airflow	<div> <div></div> <div></div> <div></div> <div></div> </div>	@daily ⌚		2023-06-23, 00:00:00 ⌚	<div> <div></div> <div></div> <div></div> <div></div> <div></div> <div></div> <div></div> <div></div> <div></div> <div></div> </div>	<div>▶</div> <div>⏻</div>	...
<div>dataset_produces_2</div> <div>dataset-scheduled produces</div>	airflow	<div> <div></div> <div></div> <div></div> <div></div> </div>	None ⌚			<div> <div></div> <div></div> <div></div> <div></div> <div></div> <div></div> <div></div> <div></div> <div></div> <div></div> </div>	<div>▶</div> <div>⏻</div>	...
<div>example_bash_operator</div> <div>example example2</div>	airflow	<div> <div></div> <div></div> <div></div> <div></div> </div>	@*-*-* ⌚		2023-06-23, 00:00:00 ⌚	<div> <div></div> <div></div> <div></div> <div></div> <div></div> <div></div> <div></div> <div></div> <div></div> <div></div> </div>	<div>▶</div> <div>⏻</div>	...

Task Instance: hello_task
at: 2023-06-24, 21:34:24 UTC



- Instance Details
- Rendered
- Log
- XCom
- List Instances, all runs
- Filter Upstream

Download Log (by attempts):

1

Task Actions

Ignore All Deps

Ignore Task State

Ignore Task Deps

Run

Past

Future

Upstream

Downstream

Recursive

Failed

Clear

Past

Future

Upstream

Downstream

Mark Failed

Past

Future

Upstream

Downstream

Mark Success

Close