# Model Experimentation by Parias

## Screenshot of all the experiments



## Screenshot of one experiment with all the artifacts visible

# Screenshot of MLflow UI after dropping features

## Screenshot of all the experiments



## Screenshot of one experiment with all the artifacts visible

# Data Pipeline

## Screenshot of successful execution Airflow DAG in graph



## Screenshot of Airflow UI grid

# Training pipeline

## Screenshot of successful execution Airflow DAG in graph



## Screenshot of Airflow UI grid

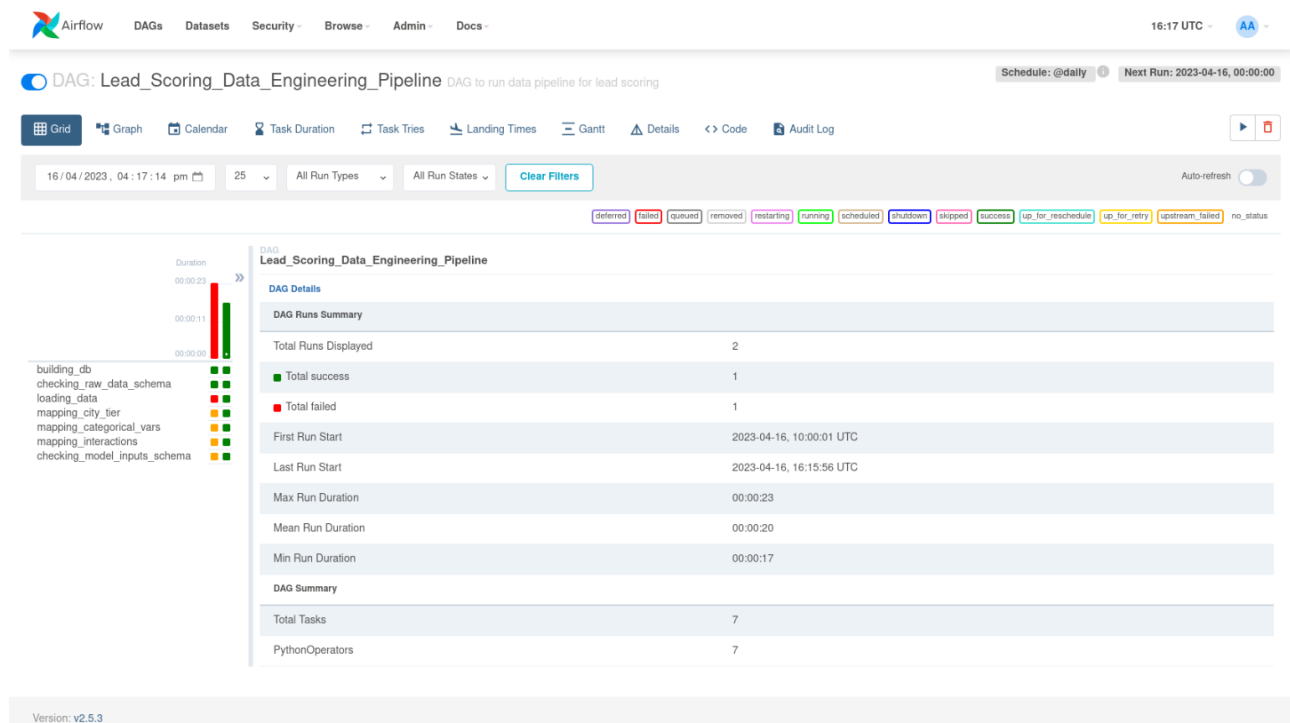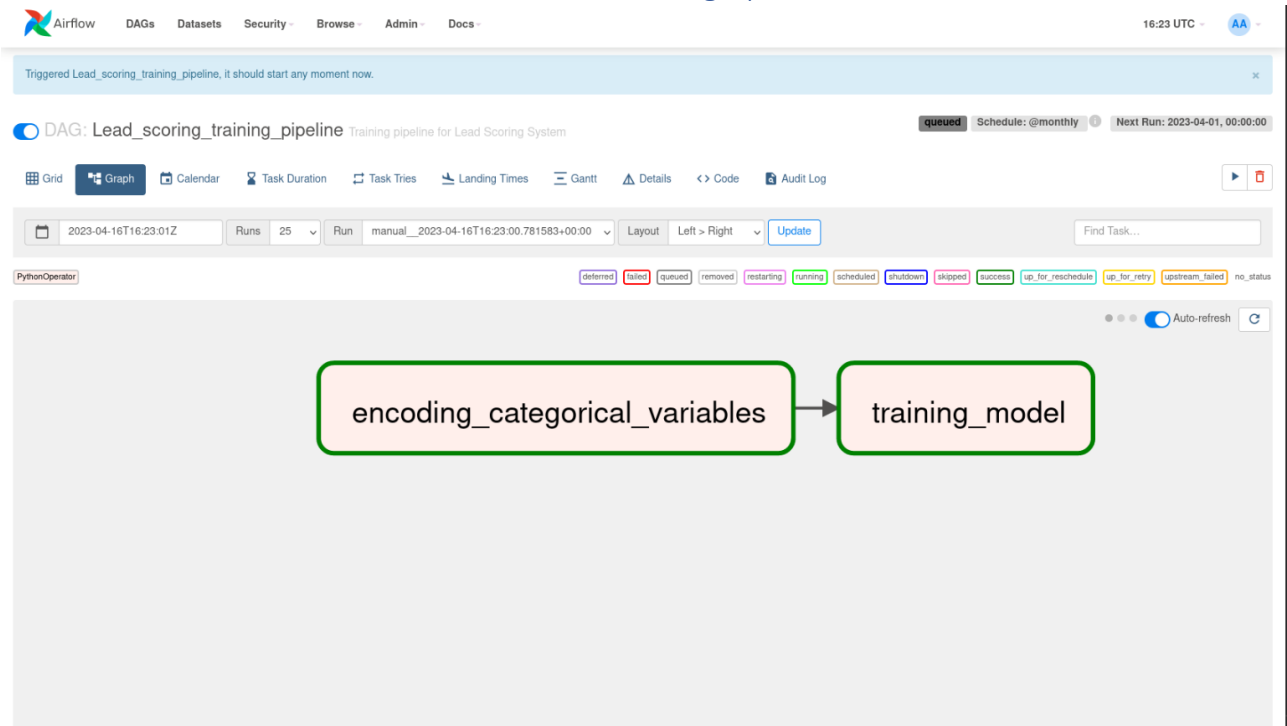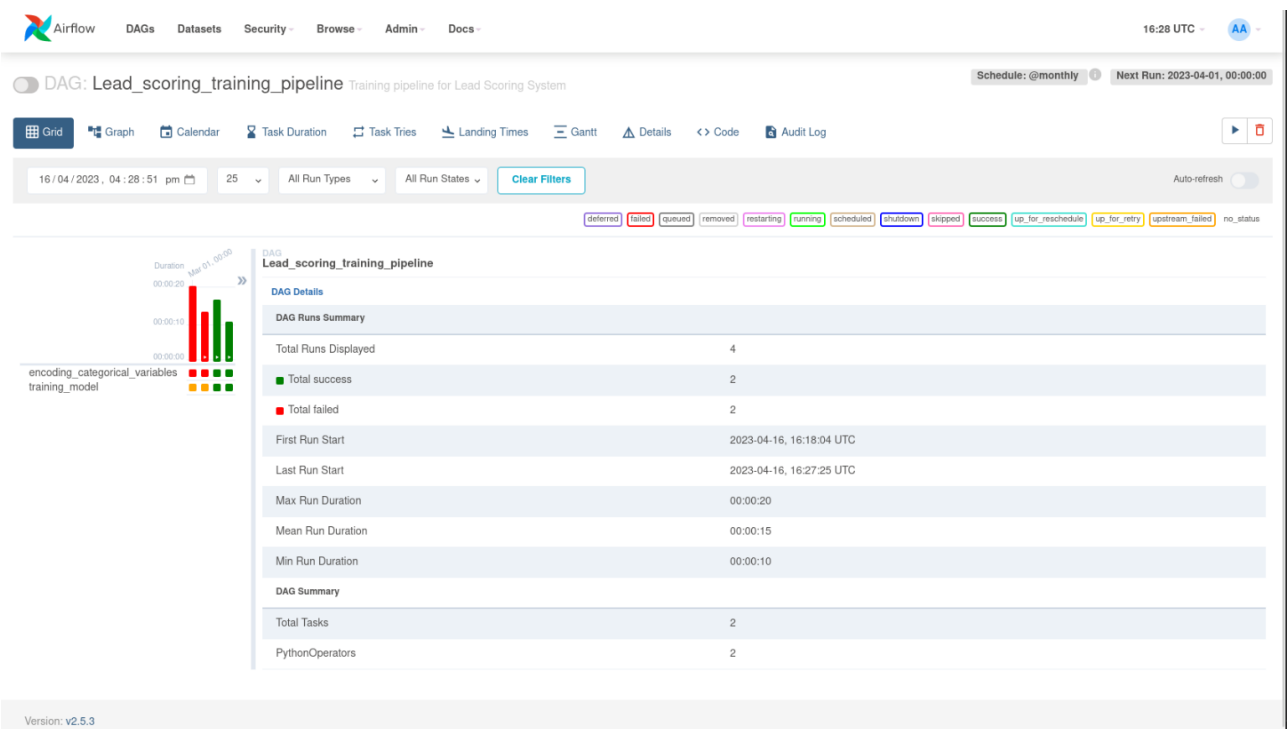# ML artifacts for training pipeline

## Screenshot of experiments with all the artifacts visible



## Screenshot of model registry with model name and stage as 'production'

# Inference Pipeline

## Screenshot of successful execution Airflow DAG in graph



## Screenshot of Airflow UI grid