

# MLOps Biomasse-Vorhersage-Pipeline

---

Dieses Repository enthält eine automatisierte Machine-Learning-Pipeline zur Vorhersage von Pflanzenbiomasse.

Die Pipeline nutzt **Dagster** für die Orchestrierung und **MLflow** für das Experiment-Tracking. Sie transformiert Rohbilder und Metadaten in ein trainiertes ResNet-Regressionsmodell.

---

## 1. Pipeline-Architektur

Die Pipeline besteht aus den folgenden Dagster-Assets:

- **raw\_dataset**  
Lädt die Metadaten (Excel) und validiert, ob die Bildpfade auf der Festplatte existieren.
  - **eda\_plots** (optional)  
Erstellt Plots zur explorativen Datenanalyse (z. B. Verteilung der Zielvariable) und speichert diese unter `figures/`.
  - **preprocessed\_data**  
Passt die Bildgröße an, normalisiert die Bilder, führt einen Train/Validation-Split durch und erstellt PyTorch DataLoaders.
  - **trained\_model**  
Trainiert ein ResNet-Modell (ResNet18 oder ResNet50).  
**Integration:** Nutzt MLflow, um Hyperparameter zu tracken, Trainingsmetriken (Loss,  $R^2$ ) pro Epoche zu loggen und das beste Modell als Artefakt zu speichern.
  - **model\_evaluation**  
Lädt das beste Modell, führt Inferenz auf dem Validierungsdatensatz durch, berechnet den MSE (Mean Squared Error) und erstellt einen „Prediction vs. Actual“-Scatter-Plot.
- 

## 2. Installationsanleitung

Voraussetzungen

- Python 3.8 oder höher
- Virtuelle Umgebung (empfohlen)

Installation

Installieren Sie die benötigten Abhängigkeiten:

```
pip install dagster dagster-webserver dagster-mlflow mlflow pandas torch
torchvision matplotlib openpyxl numpy
```

oder

```
pip install -r requirements.txt
```

Stellen Sie sicher, dass der Datensatz im Hauptverzeichnis (Root) mit der folgenden Struktur vorhanden ist:

ML0ps\_P2\_17/2025\_10\_13\_mlops\_biomass\_data/mlops\_biomass\_data/digital\_biomass\_labels.xlsx

ML0ps\_P2\_17/2025\_10\_13\_mlops\_biomass\_data/mlops\_biomass\_data/images\_med\_res/

### 3. Ausführen der Pipeline

Um die Pipeline korrekt auszuführen, müssen zwei Dienste in separaten Terminal-Fenstern gestartet werden.

#### Schritt 1: MLflow Tracking Server starten

Dieser Server übernimmt das Logging der Experimente.

```
mlflow server --port 5001
```

Zugriff auf die UI: <http://localhost:5001>

#### Schritt 2: Dagster Dev Server starten

Dieser Befehl startet den Dagster Daemon und die Weboberfläche.

```
dagster dev -f dagster_pipeline.py
```

Zugriff auf die UI: <http://localhost:3000>

### 3. Ausführung

1. Öffnen Sie die Dagster UI unter <http://localhost:3000>.
2. Navigieren Sie zum Tab Assets (oder Lineage).
3. Klicken Sie oben rechts auf "Materialize all", um die gesamte Pipeline auszuführen.

### 4. Ergebnisse

Getrackte Parameter & Metriken

Das Asset `trained_model` loggt Folgendes in MLflow:

1. **Parameter:** `epochs`, `batch_size`, `learning_rate`, `model` (z. B. `resnet18`), `optimizer`.
2. **Metriken:** `train_loss`, `val_loss`, `val_r2`.
3. **Artefakte:** Die Gewichte des besten Modells (`best_model_mlflow.pth`).

## 5. Evaluierung

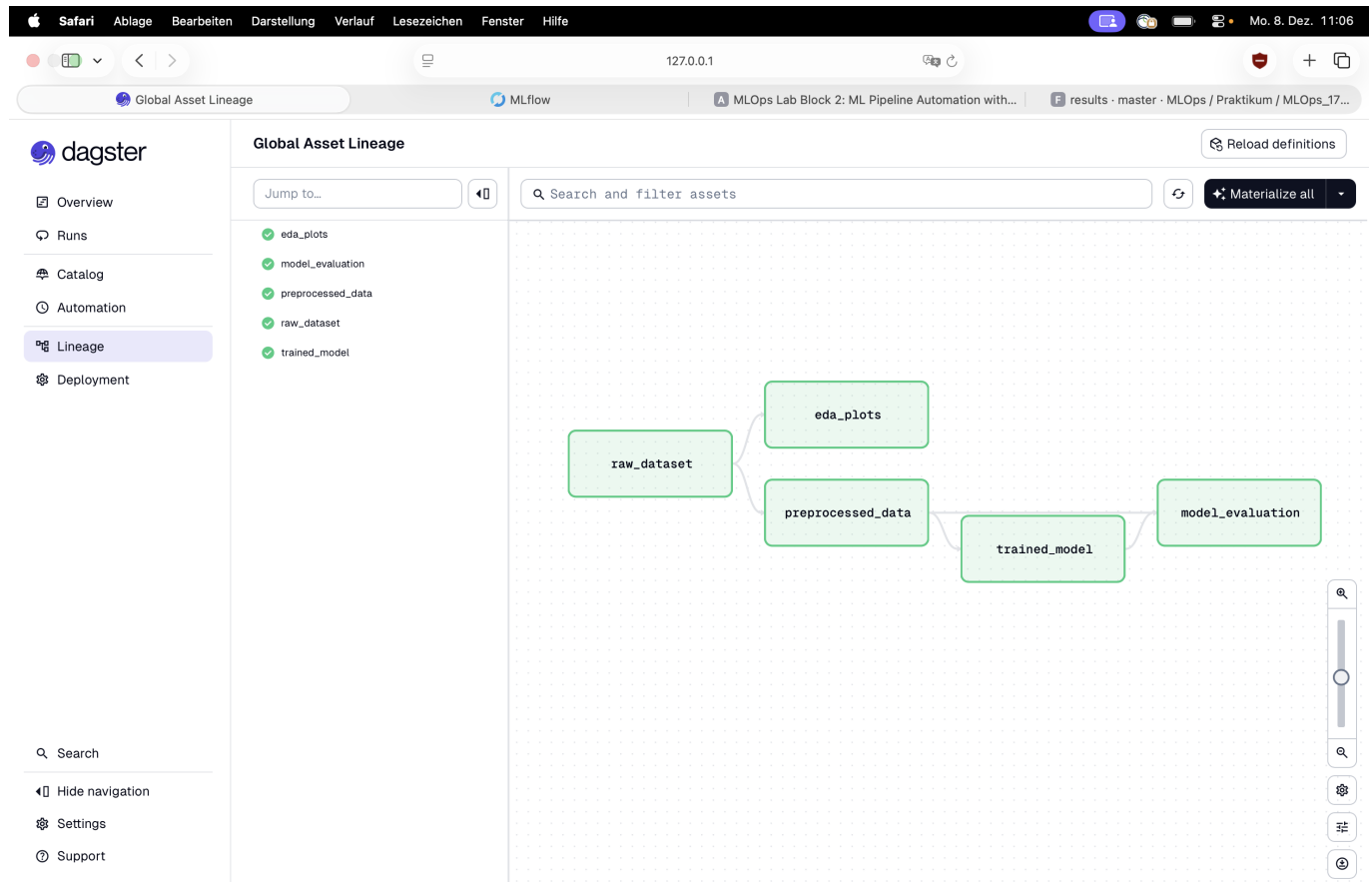
Nach einem erfolgreichen Durchlauf werden die Ergebnisse **lokal** gespeichert:

**Metriken:** results/metrics.txt (enthält den Validierungs-MSE).

**Plots:** results/pred\_vs\_actual.png und figures/target\_distribution.png.

## Screenshots

### Dagster Lineage (Erfolgreicher Durchlauf):



MLflow Experiment Tracking:

Safari

Ablage

Bearbeiten

Darstellung

Verlauf

Lesezeichen

Fenster

Hilfe

127.0.0.1

Global Asset Lineage

MLflow

MLOps Lab Block 2: ML Pipeline Automation with...

results - master - MLOps / Praktikum / MLOps\_17...

mlflow3.6.0

GitHub

Docs

plant\_biomass\_pipeline > Runs >

dagster\_resnet\_training

Overview

Model metrics

System metrics

Traces

Artifacts

Description

No description

Metrics (3)

Search metrics

Metric	Value
val_r2	0.8115066898928402
val_loss	0.014697787270921728
train_loss	0.03559899121381373

Parameters (5)

Search parameters

Parameter	Value
optimizer	Adam
learning_rate	0.001
epochs	3
model	resnet18
batch_size	32

About this run

Created at

12/08/2025, 11:00:43 AM

Created by

semih

Experiment ID

861614892808584318

Status

Finished

Run ID

ac84ae67e8fa4646b28c81fb081ee88d

Duration

2.3min

Parent run

\_\_ASSET\_JOB

Source

\_\_main\_\_.py

Logged models

Registered prompts

-

Datasets

None

Tags

Add tags

Registered models

None

Dagster Config:

Safari

Ablage

Bearbeiten

Darstellung

Verlauf

Lesezeichen

Fenster

Hilfe

127.0.0.1

MLOps Lab Block 2: ML Pipe...

screenshots - master - MLOp...

Job: \_\_ASSET\_JOB

Asset metadata and tags | D...

MLflow

Webmail :: Posteingang

dagster

Global Asset Lineage

Overview

Runs

Catalog

Automation

Lineage

Deployment

Launchpad (configure assets)

eda\_plots, model\_evaluation,

Edit tags

1 ops:

2 trained\_model:

3 config:

4 batch\_size: 32

5 epochs: 3

6 learning\_rate: 0.001

7 model\_name: resnet18

8

/\* Configure how steps are executed within a run. \*/

execution?: {

config?: ...

}

/\* Configure how loggers emit messages within a run. \*/

loggers?: {

console?: ...

}

/\* Configure runtime parameters for ops or assets. \*/

ops?: {

eda\_plots?: ...

model\_evaluation?: ...

preprocessed\_data?: ...

raw\_dataset?: ...

trained\_model?: ...

}

/\* Configure how shared resources are implemented within a run. \*/

resources?: {

/\* Built-in filesystem IO manager that stores and retrieves values using pickling. \*/

io\_manager?: ...

/\* This resource initializes an MLflow run

Use Ctrl+Space to show auto-completions inline.

}

ERRORS

No errors

CONFIG ACTIONS:

Scaffold missing config

No missing config

Scaffold all default config

All defaults expanded

Remove extra config

No extra config to remove

RUNTIME

execution

loggers

io\_manager

mlflow

ASSETS (OPS)

eda\_plots

model\_evaluation

preprocessed\_data

raw\_dataset

trained\_model

Materialize

Dagster Markdown Plot: Zum einsehen des Plots muss man auf das asset model\_evaluation klicken und unter evaluation\_plot auf "Show Markdown"

4 / 5

