# Assignment # 4

Session Fall 2025 – BSAI
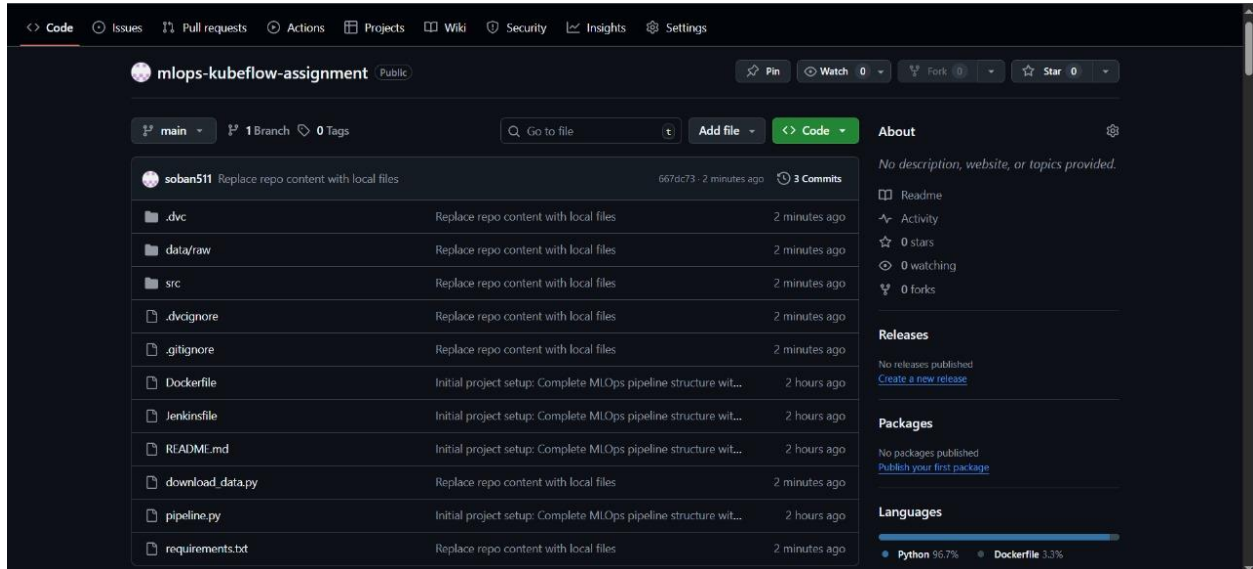
**Submitted by:**
Soban Athar 22i-0535

**Submitted to:**
Dr. Mateen Yaqoob

Department of Artificial Intelligence

National University of Computer and Emerging sciences,

FAST

# Task 1:

## Repository File Structure:



## Command "dvc status":

```
(venv) PS C:\Users\LENOVO T470\Documents\University\Semester 7\MLOPs\Assignments\Assignment4\mlops-kubeflow-assignment> dvc status
Data and pipelines are up to date.
(venv) PS C:\Users\LENOVO T470\Documents\University\Semester 7\MLOPs\Assignments\Assignment4\mlops-kubeflow-assignment>
```

## Command "dvc push":

```
(venv) PS C:\Users\LENOVO T470\Documents\University\Semester 7\MLOPs\Assignments\Assignment4\mlops-kubeflow-assignment> dvc statu
s
There are no data or pipelines tracked in this project yet.
See <https://dvc.org/doc/start> to get started!
(venv) PS C:\Users\LENOVO T470\Documents\University\Semester 7\MLOPs\Assignments\Assignment4\mlops-kubeflow-assignment> mkdir -p
../dvc-storage


    Directory: C:\Users\LENOVO T470\Documents\University\Semester 7\MLOPs\Assignments\Assignment4


Mode                 LastWriteTime         Length Name
----                 -------------         ------ ----
d-----         11/29/2025   6:15 PM                dvc-storage


(venv) PS C:\Users\LENOVO T470\Documents\University\Semester 7\MLOPs\Assignments\Assignment4\mlops-kubeflow-assignment> dvc remot
e add -d local ../dvc-storage
Setting 'local' as a default remote.
ata/raw/raw_data.csv
100% Adding...|                                                        |1/1 [00:00,  5.46file/s]

To track the changes with git, run:

        git add 'data\raw\raw_data.csv.dvc'

To enable auto staging, run:

        dvc config core.autostage true
(venv) PS C:\Users\LENOVO T470\Documents\University\Semester 7\MLOPs\Assignments\Assignment4\mlops-kubeflow-assignment> dvc push
1 file pushed
(venv) PS C:\Users\LENOVO T470\Documents\University\Semester 7\MLOPs\Assignments\Assignment4\mlops-kubeflow-assignment>
```

## Content of requirements.txt:

```
≡ requirements.txt
1    kfp==2.0.1
2    dvc==2.58.2
3    dvc-gdrive==2.20.0
4    scikit-learn==1.3.2
5    pandas==2.1.4
6    numpy==1.24.3
7    matplotlib==3.8.2
8    seaborn==0.13.0
9    joblib==1.3.2
```

# Task 2:
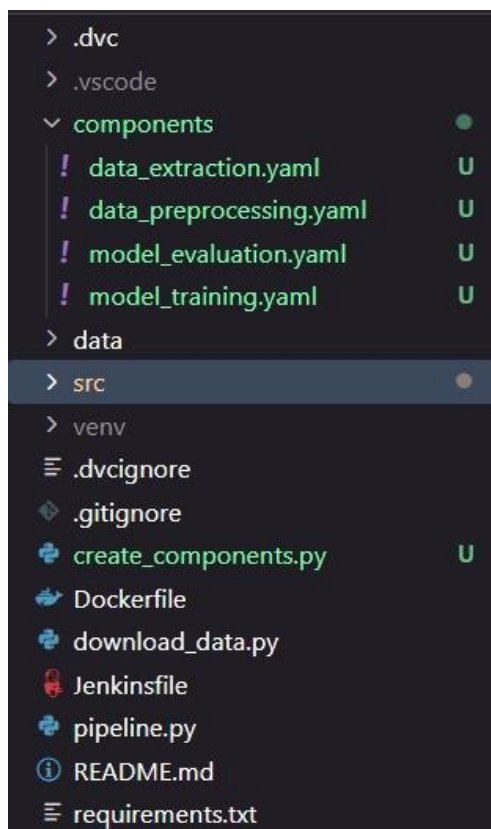
## src/pipeline_components.py file:

```python
10  @dsl.component(
11      base_image="python:3.9",
12      packages_to_install=["pandas==2.1.4", "scikit-learn==1.3.2", "numpy==1.24.3"]
13  )
14  def data_extraction(data_path: str) -> str:
15      """
16      Extracts and loads the dataset from the specified path.
17
18      Args:
19          data_path: Path to the raw data CSV file
20
21      Returns:
22          output_path: Path to the extracted data file
23
24      This component simulates data extraction. In production, you would use
25      'dvc get' or 'dvc import' to fetch versioned data from remote storage.
26      """
27      import pandas as pd
28      import os
29
30      print(f"[DATA EXTRACTION] Loading data from: {data_path}")
31
32      # Create output directory
33      output_dir = "/tmp/data"
34      os.makedirs(output_dir, exist_ok=True)
35      output_path = f"{output_dir}/extracted_data.csv"
36
37      # Load and save data
38      df = pd.read_csv(data_path)
39      df.to_csv(output_path, index=False)
40
41      print(f"[DATA EXTRACTION] Data extracted successfully")
42      print(f"[DATA EXTRACTION] Shape: {df.shape}")
43      print(f"[DATA EXTRACTION] Columns: {df.columns.tolist()}")
44      print(f"[DATA EXTRACTION] Output saved to: {output_path}")
45
46      return output_path
```

```python
@dsl.component(
    base_image="python:3.9",
    packages_to_install=["pandas==2.1.4", "scikit-learn==1.3.2", "numpy==1.24.3"]
)
def data_preprocessing(input_data_path: str) -> NamedTuple('Outputs', [
    ('train_data_path', str),
    ('test_data_path', str),
    ('feature_names', list)
]):
    """
    Preprocesses the data: handles missing values, scales features, and splits into train/test sets.

    Args:
        input_data_path: Path to the input CSV file

    Returns:
        train_data_path: Path to the processed training data
        test_data_path: Path to the processed test data
        feature_names: List of feature column names

    This component performs data cleaning, feature scaling using StandardScaler,
    and splits data into 80% training and 20% test sets.
    """
    import pandas as pd
    import numpy as np
    from sklearn.model_selection import train_test_split
    from sklearn.preprocessing import StandardScaler
    import os
    from collections import namedtuple

    print("[DATA PREPROCESSING] Starting data preprocessing...")
    print(f"[DATA PREPROCESSING] Loading data from: {input_data_path}")

    df = pd.read_csv(input_data_path)
    print(f"[DATA PREPROCESSING] Initial shape: {df.shape}")
```

**YAML files of components directory:**

# Task 3:

## Command "minicube status":



## Mlflow Pipelines UI:





## Showing the output (accuracy):

| Name | Value |
|------|-------|
| accuracy_within_10_percent | 63.725490196078425 |
| accuracy_within_20_percent | 87.25490196078431 |
| mae | 2.015223228594398 |
| mape | 10.992499453784985 |
| mean_prediction | 21.34482295993374 |
| mse | 7.772279339329266 |
| r2_score | 0.8940150227578272 |
| rmse | 2.7878807971879405 |
| std_prediction | 7.790462013463696 |

## Artifacts:

# Task 4:

## Successfull stages execution:



# Task 5:

## Repository URL:

https://github.com/soban511/mlops-kubeflow-assignment.git