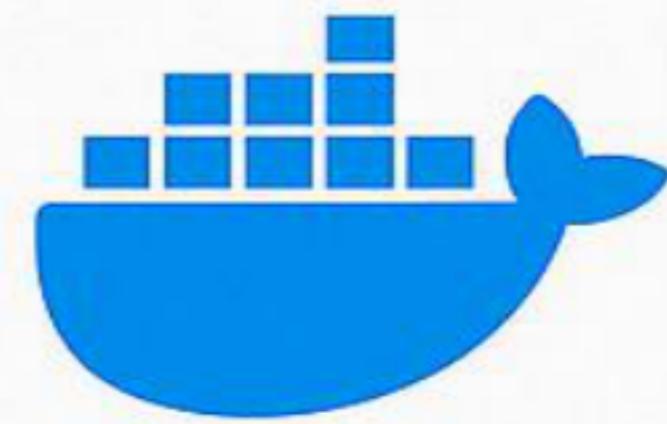
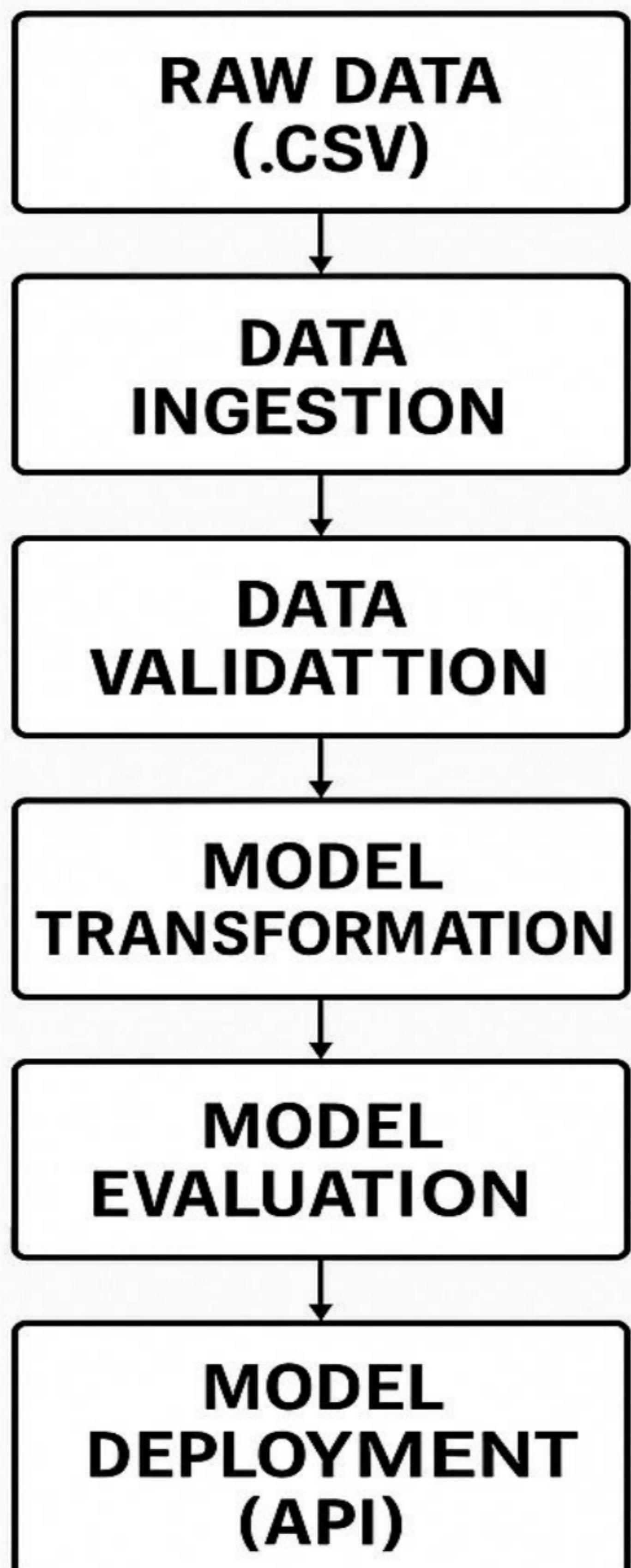


# PHISHING WEBSITE DETECTION USING A CLOUD-DEPLOYABLE ML PIPELINE



# Project Title

Network Security ML Pipeline: Phishing Website Detection and Cloud Deployment

---

## Project Overview

This end-to-end machine learning project is designed to detect phishing websites using structured features derived from URLs and network behaviour.

The solution is implemented as a binary classification pipeline, where the objective is to accurately classify a given website as either phishing (malicious) or legitimate (safe).

---

## Key Highlights of the Project

- **Problem Definition:** A supervised learning problem, specifically binary classification, aimed at identifying phishing websites based on feature-engineered data.
- **Data Source:** Structured CSV dataset containing a wide range of features such as:
  - URL characteristics (e.g., length, presence of special characters)
  - Security indicators (e.g., HTTPS, SSL status)
  - Network behavior flags
- **Architecture:** Built using a **modular pipeline structure**, split across major components:
  - Data Ingestion
  - Data Validation
  - Data Transformation
  - Model Training
  - Model Evaluation
  - Inference and Deployment
- **Robust Engineering:**
  - Custom **exception handling** and **logging mechanisms** for debuggability
  - Use of **dataclasses** to manage artifacts and configuration consistently
- **Cloud Integration:**
  - Data storage and retrieval via **MongoDB Atlas**
  - Model artifacts and metrics synced to **AWS S3**
- **Model Optimization:**
  - Hyperparameter tuning to improve **precision, recall, F1-score**
  - Tracking experiments with **MLflow**
- **Deployment Stack:**
  - REST API using **FastAPI**
  - Containerized using **Docker**
  - Deployed to **AWS EC2** for real-time predictions
- **CI/CD Automation:**
  - **GitHub Actions** automate Docker builds and push to **AWS ECR**
  - Ensures reproducible, one-click deployment pipeline

## Initial Project Structure Created in VS Code

NETWORK_SECURITY_PROJECT/	Root directory of the project
.env	Stores sensitive secrets like database passwords or API keys
.gitignore	Tells Git which files/folders to ignore (like .env, venv, __pycache__)
app.py	FastAPI app for serving predictions through API
Dockerfile	Contains Docker instructions to package the entire app
main.py	Entry point for running the entire ML training pipeline
push_data.py	Uploads raw CSV data into MongoDB
README.md	Project summary and setup instructions
requirements.txt	List of required Python libraries
setup.py	Makes the code inside `Network_security` installable as a package
test_Mongo_DB.py	Script to test MongoDB connection
venv	Python virtual environment (auto-created, not tracked in Git)
data_schema/	Contains schema definition (`schema.yaml`) to validate input CSV columns
- schema.yaml	
Network_security/	Main Python package with all logic
- __init__.py	
cloud/	Code for syncing data with AWS S3
- __init__.py	
- S3_syncer.py	
components/	All core ML pipeline steps: data ingestion, validation, etc.
- __init__.py	
- data_ingestion.py	
- data_transformation.py	
- data_validation.py	
- model_trainer.py	
constants/	Stores fixed configuration values like file paths and column names
- __init__.py	
- training_pipeline/	
__init__.py (The file contains all the constant values)	(The file contains all the constant values)
entity/	Data classes for passing configuration and output between pipeline steps
- __init__.py	
- artifact_entity.py	
- config_entity.py	
exception/	Custom exception handling for better error tracking
- __init__.py	
- exception.py	
logging/	Centralized logging setup for debugging and monitoring
- __init__.py	
- logger.py	
pipeline/	Scripts to run full training or batch prediction pipelines
- __init__.py	
- batch_prediction.py	
- training_pipeline.py	
utils/	
- __init__.py	
- main_utils/	Common utility functions used across the pipeline
__init__.py	
utils.py	
- ml_utils/	ML-specific helpers like metrics and model creation
__init__.py	
metric/	
__init__.py	
classification_report.py	
model/	
__init__.py	
estimator.py	
Network_security_data/	
- phishingData.csv	Stores the raw dataset (`phishingData.csv`) used by the pipeline

Each folder contains an `__init__.py` file so it can be used as a Python package. This modular setup ensures clear separation of concerns and allows for easy scaling and maintenance.

## Project Packaging with setup.py

```
setup.py > ...
1  from setuptools import find_packages, setup → It configures and installs project as a Python package.
2  from typing import List
3
4  def get_requirements ()->List[str]:
5      req_list:List[str]=[] → List[str] is a type hint, not part of code execution. Improves readability and static checking.
6      try:
7          with open('requirements.txt', 'r') as file:
8              lines=file.readlines()
9              for line in lines:
10                  req=line.strip()
11                  if req!='' and req!="-e .": → -e . is useful during development, but it should not be passed to install requires,
12                      req_list.append(req) because it tells pip to install the current project — which is already being installed.
13      except FileNotFoundError:
14          print("requirements.txt not found")
15      return req_list
16
17  setup(
18      name="Network_Security",
19      version="0.0.1",
20      author="kunal Shukla",
21      author_email='kpshukla3@gmail.com',
22      packages=find_packages(),
23      install_requires=get_requirements() )
24  )
```

The setup() function defines metadata (name, version, author), auto-discovers packages, and sets installation dependencies using the helper function.

## Logging and Exception Handling

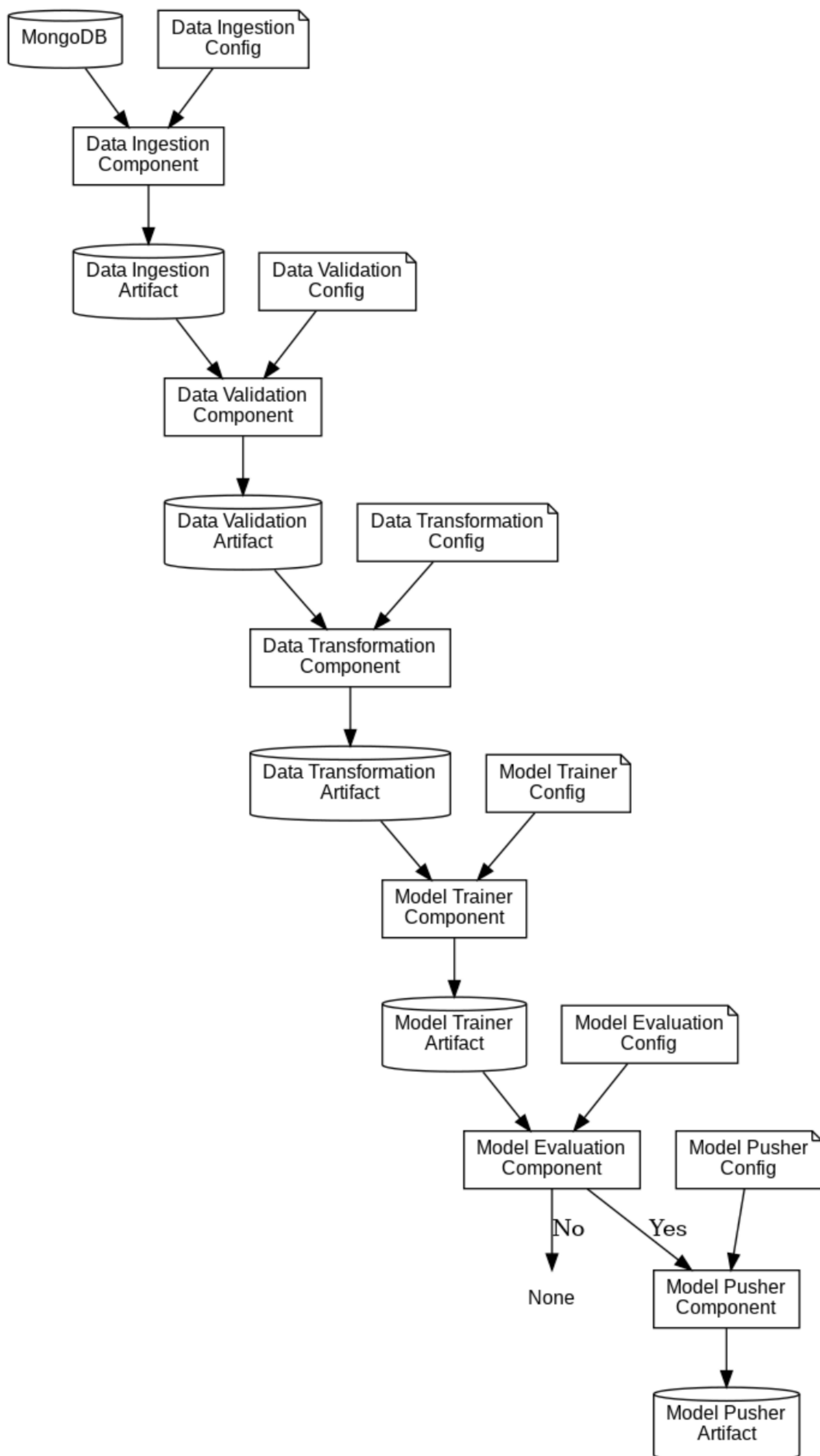
```
Network_security > logging > logger.py > ...
1  import logging
2  import os
3  from datetime import datetime
4
5  LOG_FILE=f"{datetime.now().strftime('%m_%d_%Y_%H_%M_%S')}.log" → 04_16_2025_12_35_10.log
6
7  logs_path=os.path.join(os.getcwd(),"logs",LOG_FILE)
8  os.makedirs(logs_path,exist_ok=True) → Network_Security_Project/logs/04_16_2025_12_35_10.log
9
10 LOG_FILE_PATH=os.path.join(logs_path,LOG_FILE)
11
12 logging.basicConfig(
13     filename=LOG_FILE_PATH,
14     format=" [%(asctime)s] %(lineno)d %(name)s - %(levelname)s - %(message)s",
15     level=logging.INFO,
16 )
```

Defines where logs go, what format to use, and what level of messages to record

```
Network_security > exception > exception.py > ...
1  import sys → To access system-level info
2  from Network_security.logging import logger
3
4  class NetworkSecurityException(Exception): → Inheriting from Python's built-in Exception class makes it a valid custom error that can be raised and handled like built-in exceptions.
5      def __init__(self,error_message,error_details:sys):
6          self.error_message = error_message
7          _,_,exc_tb = error_details.exc_info()
8
9          self.lineno=exc_tb.tb_lineno
10         self.file_name=exc_tb.tb_frame.f_code.co_filename
11
12     def __str__(self):
13         return "Error occurred in python script name [{0}] line number [{1}] error message [{2}].format(
14             self.file_name, self.lineno, str(self.error_message))
```

error\_details.exc\_info() returns a tuple of 3 items: (exception\_type, exception\_value, traceback\_object)  
`_ _, exc_tb` is used to retrieve only the third element

## End-to-end machine learning pipeline flowchart



## Project Dependencies with requirements.txt

```
requirements.txt
1 python-dotenv → Loads environment variables from .env file, used to manage sensitive credentials like database URIs and AWS
2 pandas
3 numpy
4 pymongo → Python client for MongoDB to connect, read, and write to the MongoDB Atlas database
5 certifi → Provides Mozilla's trusted CA certificates, used to securely connect to remote MongoDB
6 pymongo[srv]
7 scikit-learn
8 mlflow
9 pyyaml
10 fastapi
11 uvicorn
12 python-multipart → Required by FastAPI to handle file uploads in HTTP requests (used in batch prediction)
13 dill
14 -e .
15 -e .

Used to save and load complex Python objects like lambda functions or custom classes — works better than pickle in many ML use cases.
```

## Why MongoDB ?

Instead of reading data from local CSV files, we use MongoDB to store and retrieve data in a centralized, scalable, and flexible way. This makes the pipeline cloud-ready and suitable for real-time access and updates, unlike static CSVs.

## MongoDB Setup and URI Configuration

This project uses MongoDB Atlas, a cloud-based NoSQL database, for storing and retrieving structured data during ingestion and prediction. The MongoDB connection string (URI) is created in MongoDB Atlas and securely loaded from .env file using python-dotenv.

## test\_Mongo\_DB.py – Testing MongoDB Connectivity Using URI

```
.env
1 Mongo_DB_URL= generated via MongoDB Atlas
```

```
test_Mongo_DB.py > ...
1 from pymongo.mongo_client import MongoClient
2 from dotenv import load_dotenv
3 import os
4
5 # Load environment variables
6 load_dotenv()
7
8 # Get MongoDB URI from environment
9 uri = os.getenv("MONGO_DB_URL")
10 # Create a new client and connect to the server
11 client = MongoClient(uri)
12 # Send a ping to confirm a successful connection
13 try:
14     client.admin.command('ping')
15     print("Pinged your deployment. You successfully connected to MongoDB!")
16 except Exception as e:
17     print(e)
```

## push\_data.py – Converting CSV to JSON and Inserting into MongoDB

```

push_data.py > NetworkDataExtract > csv_to_json_convertor
1 import os
2 import sys
3 import json
4
5 from dotenv import load_dotenv
6 load_dotenv()
7 MONGO_DB_URL=os.getenv("MONGO_DB_URL")
8
9 import certifi
10 ca=certifi.where() ] certifi provides trusted SSL certificates
11 certifi.where() returns the path to those certificates, stored in ca, for secure MongoDB connection.
12
13 import pandas as pd
14 import numpy as np
15 import pymongo
16
17 from Network_security.exception.exception import NetworkSecurityException
18 from Network_security.logging import logger
19
20 class NetworkDataExtract():
21     def __init__(self):
22         try:
23             pass
24         except Exception as e:
25             raise NetworkSecurityException(e,sys)
26
27     def csv_to_json_convertor(self,file_path):
28         try:
29             data=pd.read_csv(file_path)
30             data.reset_index(drop=True,inplace=True)
31             records=list(json.loads(data.T.to_json()).values())
32             return records
33         except Exception as e:
34             raise NetworkSecurityException(e,sys)

```



## Phishing Website Detection – End-to-End ML Pipeline Project

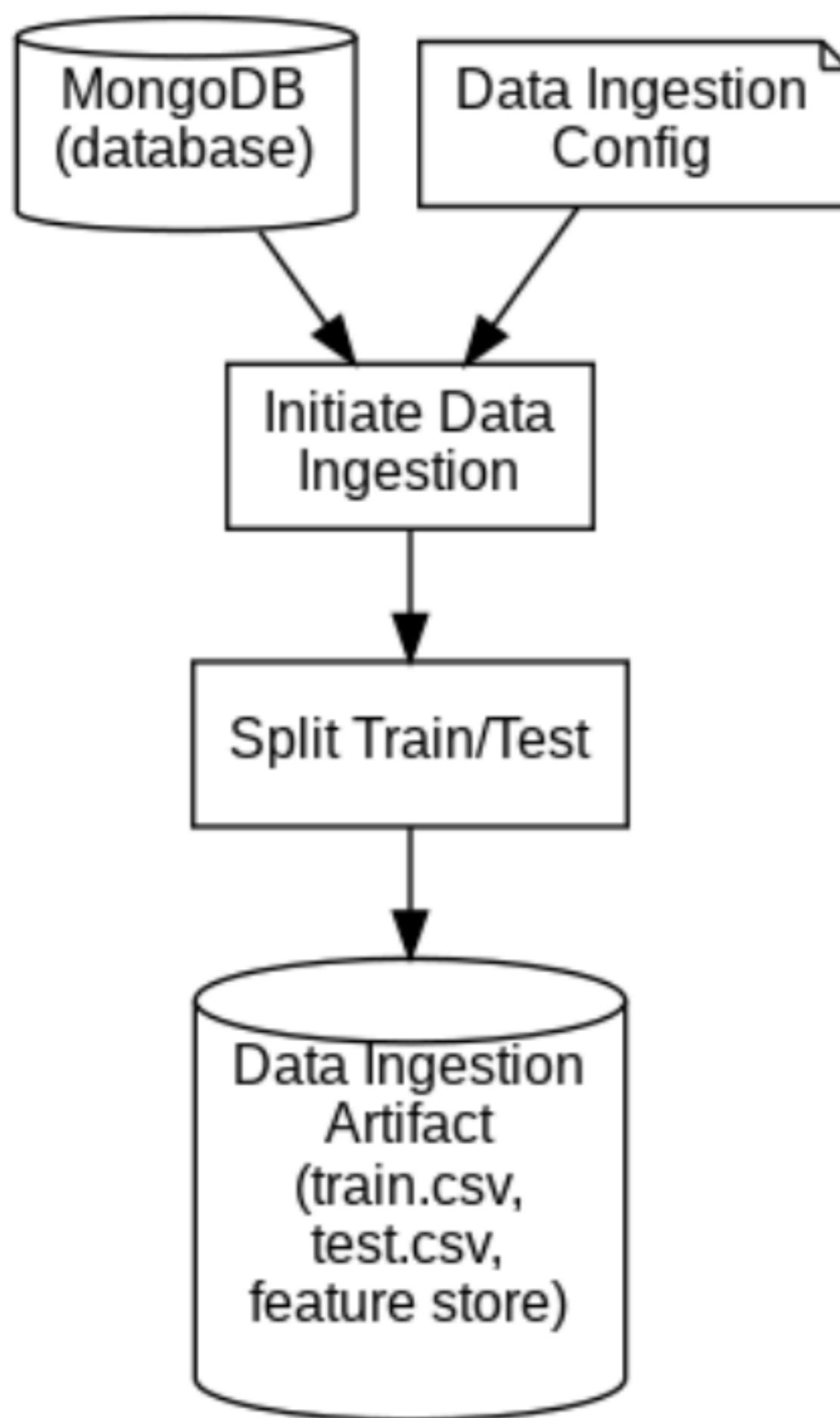
```
35     def insert_data_mongodb(self,records,database,collection):
36         try:
37             self.database=database
38             self.collection=collection
39             self.records=records
40
41             self.mongo_client=pymongo.MongoClient(MONGO_DB_URL) → Client acts as the gateway to interact with cloud database
42             self.database=self.mongo_client[self.database]
43
44             self.collection=self.database[self.collection] → Accesses the MongoDB database, which is like a schema
45             self.collection.insert_many(self.records) → that holds multiple collections.
46             return(len(self.records)) → Accesses the collection inside that database
47
48         except Exception as e:
49             raise NetworkSecurityException(e,sys)
```

Inserts all records into the MongoDB collection using batch insertion

```
50
51     if __name__=='__main__':
52         FILE_PATH="Network_security_data\phisingData.csv" → Sets the path to the input CSV file
53         DATABASE="KUNAL_AI"
54         Collection="Networkdata" ] → Defines the target MongoDB database and collection
55         networkobj=NetworkDataExtract()
56         records=networkobj.csv_to_json_convertor(file_path=FILE_PATH)
57         print(records)
58         no_of_records=networkobj.insert_data_mongodb(records,DATABASE,Collection)
59         print(no_of_records)
```

## **Data Ingestion Flow:**

### **From MongoDB to Train/Test Split**



### **training\_pipeline/init.py – Constants for Data Ingestion Configuration**

Constants are fixed values like directory names, filenames, collection names, or split ratios that do not change during execution. In this project, they are placed inside the `training_pipeline` module because they directly relate to the training process—keeping configuration centralized, organized, and easy to manage or modify later.

```
Network_security > constants > training_pipeline > _init_.py > ...
1  import os
2  import sys
3  import numpy as np
4  import pandas as pd
5
6  TARGET_COLUMN="Result"
7  PIPELINE_NAME:str="NetworkSecurity"
8  ARTIFACT_DIR:str="Artifacts"
9  FILE_NAME:str="phisingData.csv"
10
11 TRAIN_FILE_NAME:str="train.csv"
12 TEST_FILE_NAME:str="test.csv"
13
14 DATA_INGESTION_COLLECTION_NAME:str="Networkdata"
15 DATA_INGESTION_DATABASE_NAME:str="KUNAL_AI"
16 DATA_INGESTION_DIR_NAME:str="data_ingestion"
17 DATA_INGESTION_FEATURE_STORE_DIR:str="feature_store"
18 DATA_INGESTION_INGESTED_DIR:str="ingested"
19 DATA_INGESTION_TRAIN_TEST_SPLIT_RATIO:float=0.2
```

## data\_ingestion\_config.py- Setting up paths & parameters for Data Ingestion

```
Network_security > entity > config_entity.py > DataIngestionConfig > __init__
1   from datetime import datetime
2   import os
3
4   from Network_security.constants import training_pipeline → Imports predefined constants used for configuring the training
5
6   print(training_pipeline.PIPELINE_NAME)
7   print(training_pipeline.ARTIFACT_DIR)
8
9   class TrainingPipelineConfig:
10      def __init__(self,timestamp=datetime.now()):
11          timestamp=timestamp.strftime("%m_%d_%Y_%H_%M_%S")
12          self.pipeline_name=training_pipeline.PIPELINE_NAME } self.pipeline_name= NetworkSecurity
13          self.artifact_name=training_pipeline.ARTIFACT_DIR } self.artifact_name=Artifacts
14          self.artifact_dir=os.path.join(self.artifact_name,timestamp) → self.artifact_dir → Artifacts/04_17_2025_13_45_20
15          self.model_dir=os.path.join("final_model")
16          self.timestamp:str=timestamp
```

```
19  class DataIngestionConfig:
20      def __init__(self,training_pipeline_config:TrainingPipelineConfig): Type hint indicates that an object of TrainingPipelineConfig
21          self.data_ingestion_dir:str=os.path.join( must be passed during initialization
22              training_pipeline_config.artifact_dir,training_pipeline.DATA_INGESTION_DIR_NAME
23          )
24          self.feature_store_file_path:str=os.path.join(
25              self.data_ingestion_dir,training_pipeline.DATA_INGESTION_FEATURE_STORE_DIR,training_pipeline.FILE_NAME
26          )
27          self.training_file_path:str=os.path.join(
28              self.data_ingestion_dir,training_pipeline.DATA_INGESTION_INGESTED_DIR,training_pipeline.TRAIN_FILE_NAME
29          )
30          self.testing_file_path:str=os.path.join(
31              self.data_ingestion_dir,training_pipeline.DATA_INGESTION_INGESTED_DIR,training_pipeline.TEST_FILE_NAME
32          )
33          self.train_test_split_ratio:float=training_pipeline.DATA_INGESTION_TRAIN_TEST_SPLIT_RATIO
34          self.collection_name:str=training_pipeline.DATA_INGESTION_COLLECTION_NAME
35          self.database_name:str=training_pipeline.DATA_INGESTION_DATABASE_NAME
```

self.data\_ingestion\_dir → Artifacts/04\_17\_2025\_13\_45\_20/data\_ingestion  
 self.feature\_store\_file\_path → Artifacts/04\_17\_2025\_13\_45\_20/data\_ingestion/feature\_store/phisingData.csv  
 self.training\_file\_path → Artifacts/04\_17\_2025\_13\_45\_20/data\_ingestion/ingested/train.csv  
 self.testing\_file\_path → Artifacts/04\_17\_2025\_13\_45\_20/data\_ingestion/ingested/test.csv

# Phishing Website Detection – End-to-End ML Pipeline Project

## data\_ingestion.py – Reading Data from MongoDB and Splitting into Train/Test

This module handles the complete data ingestion process. It pulls raw data from a MongoDB collection, stores it in a feature store (as CSV), and performs a train-test split. It returns a DataIngestionArtifact for downstream pipeline components.

```
Network_security > components > data_ingestion.py > DataIngestion > export_collection_as_dataframe
1  import pymongo.mongo_client
2  from Network_security.exception.exception import NetworkSecurityException
3  from Network_security.logging.logger import logging
4
5  from Network_security.entity.config_entity import DataIngestionConfig
6  from Network_security.entity.artifact_entity import DataIngestionArtifact
7
8  import os
9  import sys
10 import numpy as np
11 import pandas as pd
12 import pymongo
13 from typing import List
14 from sklearn.model_selection import train_test_split
15
16 from dotenv import load_dotenv
17 load_dotenv()
18 Mongo_DB_URL=os.getenv("Mongo_DB_URL")
19
20 class DataIngestion:
21     def __init__(self,data_ingestion_config:DataIngestionConfig):
22         try:
23             self.data_ingestion_config=data_ingestion_config
24         except Exception as e:
25             raise NetworkSecurityException(e,sys)

27     def export_collection_as_dataframe(self):
28         try:
29             database_name=self.data_ingestion_config.database_name
30             collection_name=self.data_ingestion_config.collection_name
31             self.mongo_client=pymongo.MongoClient(Mongo_DB_URL)
32             collection=self.mongo_client[database_name][collection_name]
33             df=pd.DataFrame(list(collection.find()))
34             if "_id" in df.columns.to_list():
35                 df=df.drop(columns=["_id"],axis=1)
36             df.replace({"na":np.nan},inplace=True)
37             return df
38         except Exception as e:
39             raise NetworkSecurityException(e,sys)
40
41     def export_data_into_feature_store(self,dataframe:pd.DataFrame):
42         try:
43             feature_store_file_path=self.data_ingestion_config.feature_store_file_path
44             # Creating a folder
45             dir_path=os.path.dirname(feature_store_file_path)
46             os.makedirs(dir_path,exist_ok=True)
47             dataframe.to_csv(feature_store_file_path,index=False,header=True)
48             return dataframe
49         except Exception as e:
50             raise NetworkSecurityException(e,sys)
```

Required for initializing data ingestion: config provides input settings, artifact stores output file paths

Retrieves database and collection names from the ingestion config object

Fetches all documents from the collection and converts them into DataFrame

Removes the MongoDB default \_id column, which is not needed for model

Replaces any "na" string entries with proper NaN values for compatibility

Artifacts/04\_17\_2025\_13\_45\_20/data\_ingestion/feature\_store/phisingData.csv

We convert the MongoDB data to CSV so it can be stored locally in the feature store and reused consistently across the pipeline, without repeatedly querying the database.

## Phishing Website Detection – End-to-End ML Pipeline Project

```
52     def split_data_as_train_test(self,dataframe:pd.DataFrame):
53         try:
54             train_set,test_set=train_test_split(
55                 dataframe,test_size=self.data_ingestion_config.train_test_split_ratio)
56         )
57         logging.info("Performed train test split on DataFrame")
58         logging.info(
59             "Exited split_data_as_train_test method of Data_Ingestion class"
60         )
61         dir_path=os.path.dirname(self.data_ingestion_config.training_file_path)
62         os.makedirs(dir_path,exist_ok=True)
63         logging.info(f"Exporting train and test file path,")
64
65         train_set.to_csv(
66             self.data_ingestion_config.training_file_path,index=False,header=True
67         )
68         test_set.to_csv(
69             self.data_ingestion_config.testing_file_path,index=False,header=True
70         )
71         logging.info(f"Exporting train and test file path,")
72
73     except Exception as e:
74         raise NetworkSecurityException(e,sys)
```

Splits data into training and testing sets using ratio from config

Artifacts/04\_17\_2025\_13\_45\_20/data\_ingestion/ingested/ folder will be created to store train.csv and test.csv

Artifacts/04\_17\_2025\_13\_45\_20/data\_ingestion/ingested/train.csv

Artifacts/04\_17\_2025\_13\_45\_20/data\_ingestion/ingested/test.csv

```
77     def initiate_data_ingestion(self):
78         try:
79             dataframe=self.export_collection_as_dataframe()
80             dataframe=self.export_data_into_feature_store(dataframe)
81             self.split_data_as_train_test(dataframe)
82             dataingestionartifact=DataIngestionArtifact(
83                 trained_file_path=self.data_ingestion_config.training_file_path,
84                 test_file_path=self.data_ingestion_config.testing_file_path)
85             return dataingestionartifact
86
87         except Exception as e:
88             raise NetworkSecurityException(e,sys)
```

Fetching data from MongoDB into a DataFrame

Saving raw data into a local feature store CSV

Splitting data into train and test sets

Returning a DataIngestionArtifact, containing paths to train.csv and test.csv

### artifact\_entity.py – DataIngestionArtifact Definition

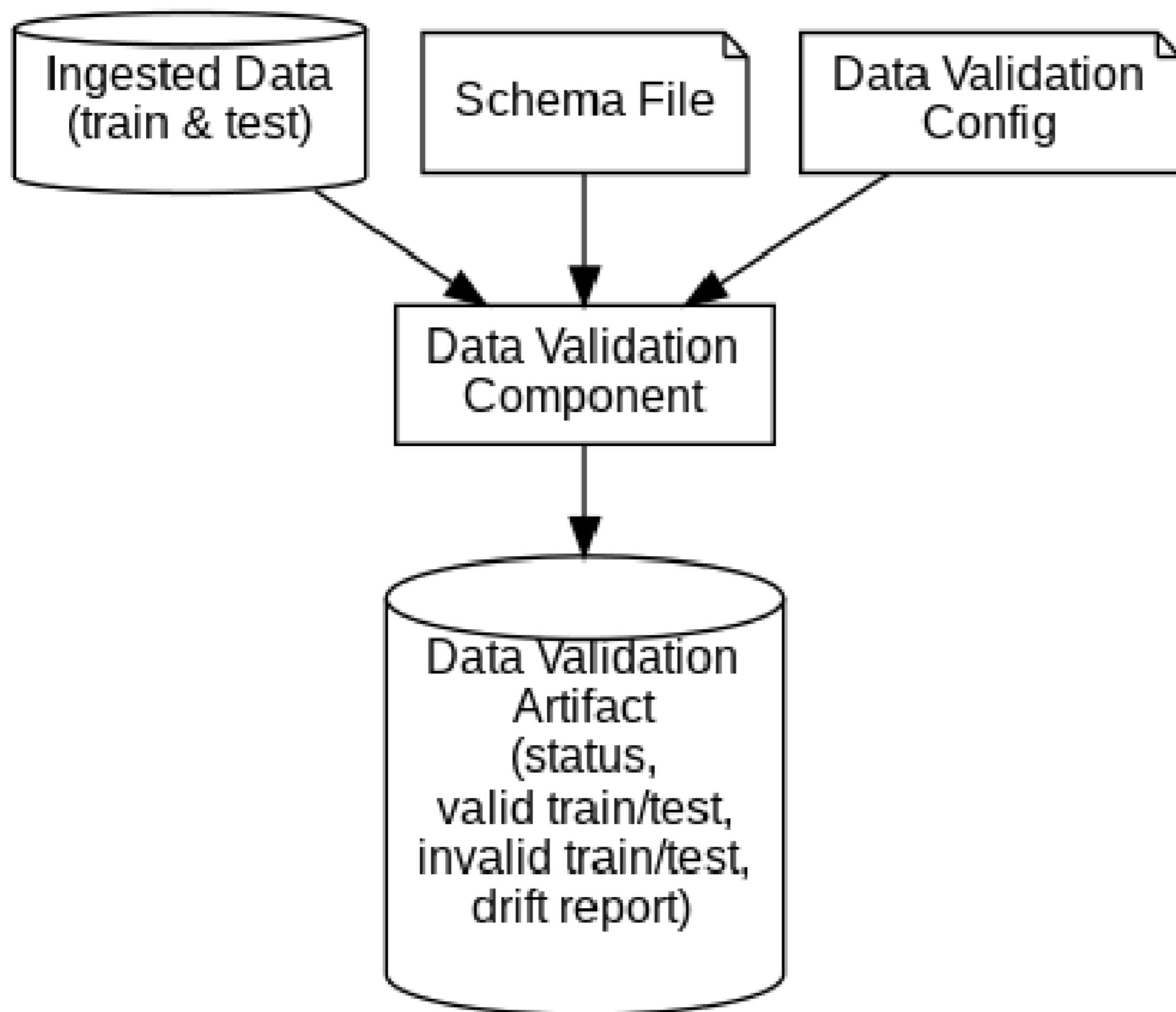
This dataclass represents the output (artifact) of the data ingestion step. It holds the file paths for both the training and testing datasets, which are used as inputs for the next stages of the ML pipeline.

```
Network_security > entity > artifact_entity.py > DataTransformationArtifact > transformed_train_file_path
1   from dataclasses import dataclass
2
3   @dataclass
4   class DataIngestionArtifact:
5       trained_file_path:str
6       test_file_path:str
```

Instead of returning artifacts directly from data ingestion, we define a structured dataclass here to store them and access them consistently across the pipeline.

## Data Validation Flow:

### From Schema Check to Drift Detection



### training\_pipeline/init.py – Constants for Data Validation Configuration

```

Network_security > constants > training_pipeline > _init_.py > ...
25 DATA_VALIDATION_DIR_NAME: str="data validation"
26 DATA_VALIDATION_VALID_DIR: str="validated"
27 DATA_VALIDATION_INVALID_DIR: str="invalid"
28 DATA_VALIDATION_DRIFT_REPORT_DIR: str="drift_report"
29 DATA_VALIDATION_DRIFT_REPORT_DIR_NAME: str="report.yaml"
30 PREPROCESSING_OBJECT_FILE_NAME="preprocessing.pkl"
  
```

### data\_validation\_config.py – defining rules and paths for Data Validation

```

Network_security > entity > config_entity.py > DataTransformationConfig > _init_
37 class DataValidationConfig:
38     def __init__(self, training_pipeline_config:TrainingPipelineConfig):
39         self.data_validation_dir:str=os.path.join(training_pipeline_config.artifact_dir,
40                                         training_pipeline.DATA_VALIDATION_DIR_NAME)
41         self.valid_data_dir:str=os.path.join(self.data_validation_dir,training_pipeline.DATA_VALIDATION_VALID_DIR)
42         self.invalid_data_dir:str=os.path.join(self.data_validation_dir,training_pipeline.DATA_VALIDATION_INVALID_DIR)
43         self.valid_train_file_path:str=os.path.join(self.valid_data_dir,training_pipeline.TRAIN_FILE_NAME)
44         self.valid_test_file_path:str=os.path.join(self.valid_data_dir,training_pipeline.TEST_FILE_NAME)
45         self.invalid_train_file_path:str=os.path.join(self.invalid_data_dir,training_pipeline.TRAIN_FILE_NAME)
46         self.invalid_test_file_path:str=os.path.join(self.invalid_data_dir,training_pipeline.TEST_FILE_NAME)
47         self.drift_report_file_path:str=os.path.join(self.data_validation_dir,
48                                         training_pipeline.DATA_VALIDATION_DRIFT_REPORT_DIR,
49                                         training_pipeline.DATA_VALIDATION_DRIFT_REPORT_DIR_NAME)

        self.data_validation_dir → Artifacts/04_17_2025_13_45_20/data validation
        self.valid_data_dir → Artifacts/04_17_2025_13_45_20/data validation/validated
        self.invalid_data_dir → Artifacts/04_17_2025_13_45_20/data validation/invalid
        self.valid_train_file_path → Artifacts/04_17_2025_13_45_20/data validation/validated/train.csv
        self.valid_test_file_path → Artifacts/04_17_2025_13_45_20/data validation/validated/test.csv
        self.invalid_train_file_path → Artifacts/04_17_2025_13_45_20/data validation/invalid/train.csv
        self.invalid_test_file_path → Artifacts/04_17_2025_13_45_20/data validation/invalid/test.csv
        self.drift_report_file_path → Artifacts/04_17_2025_13_45_20/data validation/drift_report/report.yaml
  
```

## schema.yaml – Defining expected columns and numerical fields for validation

This file lists all required columns along with their expected data types. The numerical\_columns section is added separately to help the validation component detect numerical drift. Both sections are used in data\_validation.py to ensure the structure and types match the expected format before training begins.

```
data_schema > ! schema.yaml
1   columns:
2     - having_IP_Address: int64
3     - URL_Length: int64
4     - Shortining_Service: int64
5     - having_At_Symbol: int64
6     - double_slash_redirecting: int64
7     - Prefix_Suffix: int64
8     - having_Sub_Domain: int64
9     - SSLfinal_State: int64
10    - Domain_registration_length: int64
11    - Favicon: int64
12    - port: int64
13    - HTTPS_token: int64
14    - Request_URL: int64
15    - URL_of_Anchor: int64
16    - Links_in_tags: int64
17    - SFH: int64
18    - Submitting_to_email: int64
19    - Abnormal_URL: int64
20    - Redirect: int64
21    - on_mouseover: int64
22    - RightClick: int64
```

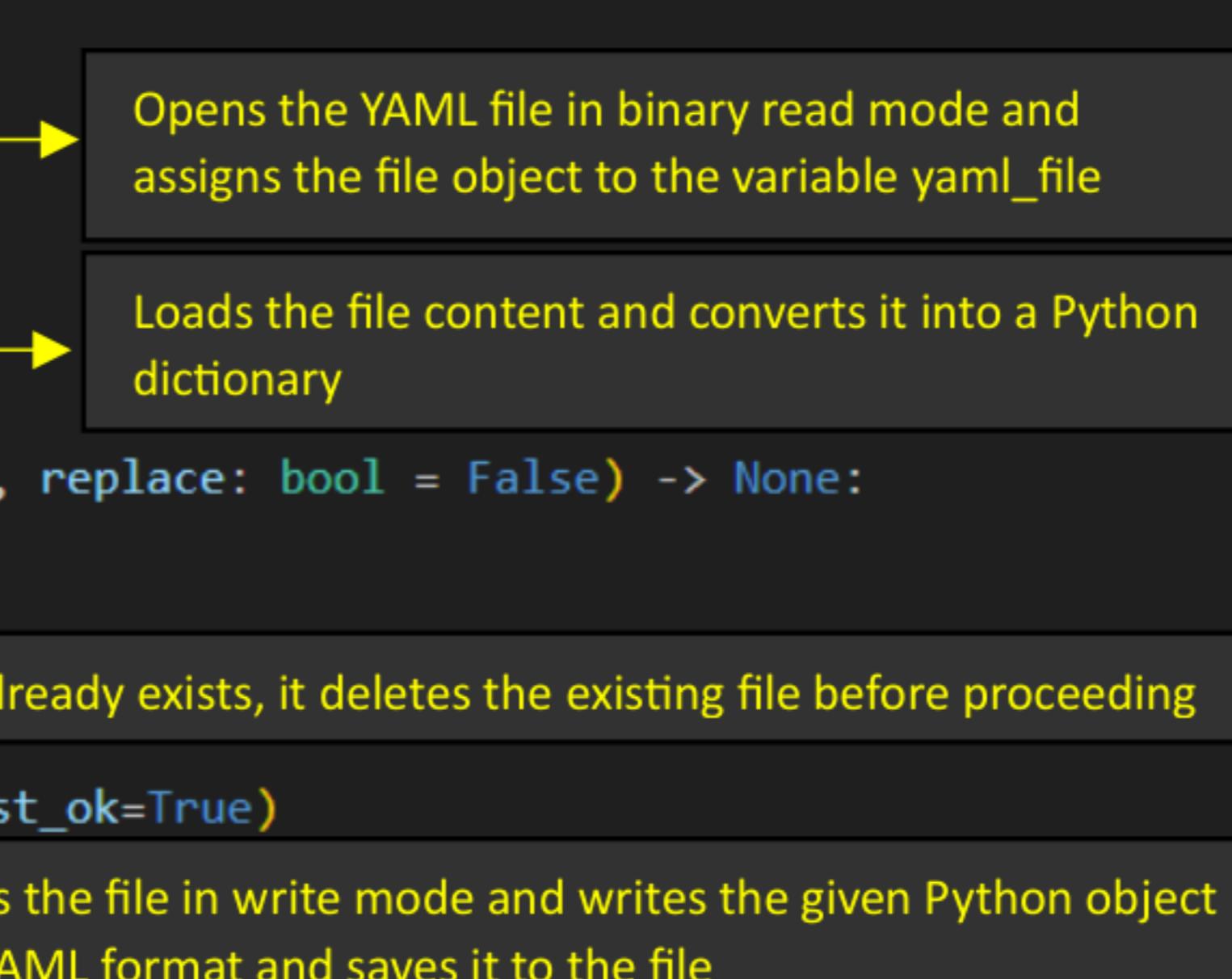
  

```
data_schema > ! schema.yaml
35  numerical_columns:
36    - having_IP_Address
37    - URL_Length
38    - Shortining_Service
39    - having_At_Symbol
40    - double_slash_redirecting
41    - Prefix_Suffix
42    - having_Sub_Domain
43    - SSLfinal_State
44    - Domain_registration_length
45    - Favicon
46    - port
47    - HTTPS_token
48    - Request_URL
49    - URL_of_Anchor
50    - Links_in_tags
51    - SFH
52    - Submitting_to_email
53    - Abnormal_URL
54    - Redirect
55    - on_mouseover
56    - RightClick
```

## utils – Supporting Functions for File Handling, Model Saving, and Metric Evaluation

The utils directory contains helper scripts that are used across multiple pipeline components. It is organized into submodules to separate general-purpose utilities and ML-specific operations.

```
Network_security > utils > main_utils > ! utils.py > ! save_numpy_array > ! dir_path
1  import yaml
2  from Network_security.exception.exception import NetworkSecurityException
3  from Network_security.logging.logger import logging
4  import os,sys
5  import numpy as np
6  import dill
7  import pickle
8  from sklearn.metrics import accuracy_score
9  from sklearn.model_selection import GridSearchCV
10 def read_yaml_file(file_path: str) -> dict:
11     try:
12         with open(file_path, "rb") as yaml_file:
13             return yaml.safe_load(yaml_file)
14     except Exception as e:
15         raise NetworkSecurityException(e, sys)
16 def write_yaml_file(file_path: str, content: object, replace: bool = False) -> None:
17     try:
18         if replace:
19             if os.path.exists(file_path):
20                 os.remove(file_path)
21             os.makedirs(os.path.dirname(file_path), exist_ok=True)
22             with open(file_path, "w") as file:
23                 yaml.dump(content, file)
24     except Exception as e:
25         raise NetworkSecurityException(e, sys)
```



Annotations for the code:

- Line 12: Opens the YAML file in binary read mode and assigns the file object to the variable `yaml_file`
- Line 13: Loads the file content and converts it into a Python dictionary
- Line 19: If file already exists, it deletes the existing file before proceeding
- Line 22: Opens the file in write mode and writes the given Python object into YAML format and saves it to the file

# Phishing Website Detection – End-to-End ML Pipeline Project

## data\_validation.py – Validating Schema and Detecting Data Drift

```
Network_security > components > data_validation.py > DataValidation > detect_database_drift
1  from Network_security.entity.artifact_entity import DataIngestionArtifact,DataValidationArtifact
2  from Network_security.entity.config_entity import DataValidationConfig
3
4  from Network_security.logging.logger import logging
5  from Network_security.exception.exception import NetworkSecurityException
6  from Network_security.constants.training_pipeline import SCHEMA_FILE_PATH
7  from scipy.stats import ks_2samp
8  import pandas as pd      ────────── Used to compare the distribution of numerical columns between train and test sets to detect data drift
9  import os, sys
10 from Network_security.utils.main_utils.utils import read_yaml_file,write_yaml_file
11
12 class DataValidation:
13     def __init__(self,data_ingestion_artifact:DataIngestionArtifact,
14                  data_validation_config:DataValidationConfig):
15         try:
16             self.data_ingestion_artifact=data_ingestion_artifact
17             self.data_validation_config=data_validation_config
18             self.schema_config=read_yaml_file(SCHEMA_FILE_PATH)
19         except Exception as e:
20             raise NetworkSecurityException(e,sys)
```

```
Network_security > components > data_validation.py > DataValidation > detect_database_drift
12  class DataValidation:
13      @staticmethod
14      def read_data(file_path)->pd.DataFrame:
15          try:
16              return pd.read_csv(file_path)
17          except Exception as e:
18              raise NetworkSecurityException(e,sys) ────────── Reads a CSV file using pd.read_csv() and returns a DataFrame; defined as a static method since it doesn't use class variables.
29  def validate_number_of_columns(self,dataframe:pd.DataFrame)->bool:
30      try:
31          number_of_columns=len(self.schema_config['columns']) ────────── Gets the expected number of columns from the schema file
32          logging.info(f"Required no of columns:{number_of_columns}")
33          logging.info(f"Data frame columns:{len(dataframe.columns)}")
34          if len(dataframe.columns)==number_of_columns:
35              return True
36          return False
37      except Exception as e:
38          raise NetworkSecurityException(e,sys) ────────── Compares the number of columns in the DataFrame with the expected count and returns True if they match.
```

```
Network_security > components > data_validation.py > DataValidation > detect_database_drift
12  class DataValidation:
40      def validate_number_of_numerical_columns(self,dataframe:pd.DataFrame)->bool:
41          try:
42              num_col_list=[]
43              for i in dataframe.columns:
44                  if dataframe[i].dtype=='int64':
45                      num_col_list.append(i) ────────── Loops through the DataFrame columns and collects all columns with int64 type into num_col_list.
47              number_of_numerical_columns=len(self.schema_config['numerical_columns']) ────────── Gets the expected count of numerical columns in the schema
48              logging.info(f"Required no of numerical columns:{number_of_numerical_columns}")
49              logging.info(f"Data frame numerical columns:{len(num_col_list)}")
50              if len(num_col_list)==number_of_numerical_columns:
51                  return True
52              return False
53          except Exception as e:
54              raise NetworkSecurityException(e,sys) ────────── Returns True if the count of int64 columns matches the number defined in the schema.
```

# Phishing Website Detection – End-to-End ML Pipeline Project

```
Network_security > components > data_validation.py > DataValidation > initiate_data_validation
12     class DataValidation:
56         def detect_database_drift(self,base_df,current_df,threshold=0.05)->0.05:
57             try:
58                 status=True
59                 report={}
60                 for column in base_df.columns:
61                     d1=base_df[column]
62                     d2=current_df[column]
63                     is_same_dist=ks_2samp(d1,d2) → Iterates over each column in the base dataset and
64                     if threshold<is_same_dist.pvalue: retrieves corresponding columns from both DataFrames.
65                         is_found=False
66                     else: → Performs a Kolmogorov-Smirnov test to compare the distribution of the column in both datasets
67                         is_found=True
68                         status=False
69                     report.update({column:{})
70                         'p_value':float(is_same_dist.pvalue),
71                         'drift_status':is_found
72                     }})
73                     drift_report_file_path=self.data_validation_config.drift_report_file_path
74                     dir_path=os.path.dirname(drift_report_file_path)
75                     os.makedirs(dir_path,exist_ok=True)
76                     write_yaml_file(file_path=drift_report_file_path,content=report) → If p-value is less than the threshold (default 0.05), drift is
77                     drift_report_file_path = Artifacts/04_17_2025_13_45_20/data validation/drift_report/report.yaml
78             except Exception as e:
79                 raise NetworkSecurityException(e,sys)
```

```
Network_security > components > data_validation.py > DataValidation > initiate_data_validation
12     class DataValidation:
85         def initiate_data_validation(self)->DataValidationArtifact:
86             try:
87                 train_file_path=self.data_ingestion_artifact.trained_file_path
88                 test_file_path=self.data_ingestion_artifact.test_file_path
89                 train_dataframe=DataValidation.read_data(train_file_path)
90                 test_dataframe=DataValidation.read_data(test_file_path)
91                 status=self.validate_number_of_columns(dataframe=train_dataframe) → Reads train and test CSV files from
92                 if not status: ingestion artifacts and loads them as
93                     error_message=f"Train dataframe does not contain all columns. \n"
94                     status=self.validate_number_of_columns(dataframe=test_dataframe) DataFrames
95                     if not status:
96                         error_message=f"Test dataframe does not contain all columns. \n"
97                         status=self.validate_number_of_numerical_columns(dataframe=train_dataframe)
98                         if not status:
99                             error_message=f"Train dataframe does not contain same numbers of numerical columns. \n"
100                            status=self.validate_number_of_numerical_columns(dataframe=test_dataframe)
101                            if not status:
102                                error_message=f"Test dataframe does not contain same numbers of numerical columns. \n"
```

```
Network_security > components > data_validation.py > DataValidation > initiate_data_validation
12     class DataValidation:
85         def initiate_data_validation(self)->DataValidationArtifact:
104             status=self.detect_database_drift(base_df=train_dataframe,current_df=test_dataframe)
105             dir_path=os.path.dirname(self.data_validation_config.valid_train_file_path)
106             os.makedirs(dir_path,exist_ok=True)
107             train_dataframe.to_csv(self.data_validation_config.valid_train_file_path,index=False,header=True)
108             test_dataframe.to_csv(self.data_validation_config.valid_test_file_path,index=False,header=True)
109             data_validation_artifact = DataValidationArtifact(
110                 validation_status=status,
111                 valid_train_file_path=self.data_ingestion_artifact.trained_file_path,
112                 valid_test_file_path=self.data_ingestion_artifact.test_file_path,
113                 invalid_train_file_path=None,
114                 invalid_test_file_path=None,
115                 drift_report_file_path=self.data_validation_config.drift_report_file_path,
116             )
117             return data_validation_artifact
118         except Exception as e:
119             raise NetworkSecurityException(e,sys)
```

## artifact\_entity.py – DataValidationArtifact Definition

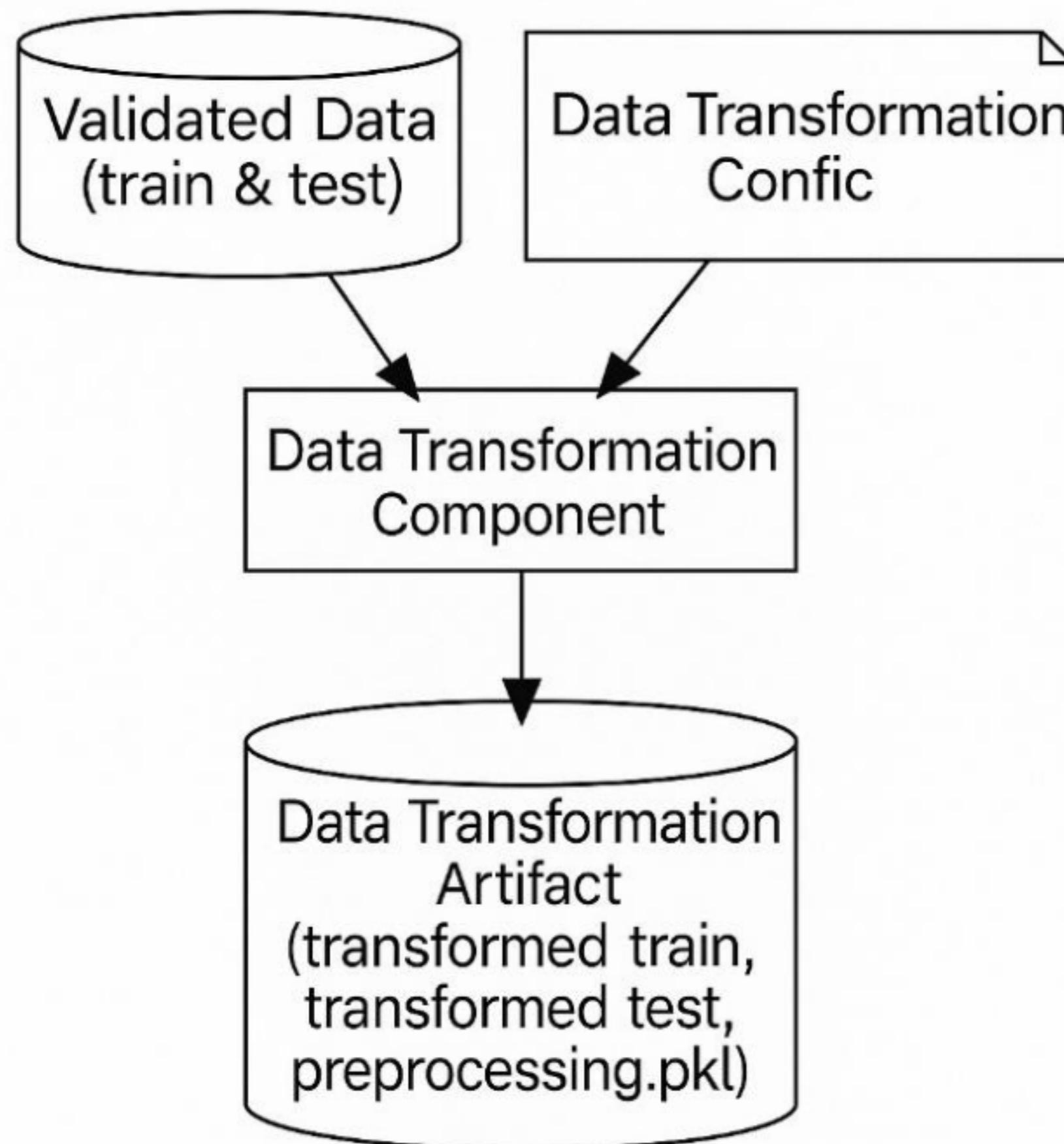
Defines the structure for storing data validation outputs like validation status, valid/invalid file paths, and drift report location. This artifact is returned at the end of the validation component.

```
Network_security > entity > artifact_entity.py > DataTransformationArtifact > transformed_train_file_path

8  @dataclass
9  class DataValidationArtifact:
10     validation_status:bool
11     valid_train_file_path:str
12     valid_test_file_path:str
13     invalid_train_file_path:str
14     invalid_test_file_path:str
15     drift_report_file_path:str
```

## Data Transformation Flow:

### Transforming and Saving Data with preprocessing.pkl



### training\_pipeline/init.py - Constants for Data Transformation Configuration

```
Network_security > constants > training_pipeline > _init_.py > ...
30 PREPROCESSING_OBJECT_FILE_NAME="preprocessing.pkl"
31
32 DATA_TRANSFORMATION_DIR_NAME:str="data_transformation"
33 DATA_TRANSFORMATION_TRANSFORMED_DATA_DIR:str="transformed"
34 DATA_TRANSFORMATION_TRANSFORMED_OBJECT_DIR:str="transformed_object"
35 # knn Imputer for nan values
36 DATA_TRANSFORMATION_IMPUTER_PARAMS: dict={}
37     "missing_values":np.nan,
38     "n_neighbors":3,
39     "weights":"uniform",
40 }
```

KNNImputer replaces missing values using the average of the 3 nearest rows based on feature similarity, preserving local data patterns. It's more contextual than mean/median imputation and works well when data isn't too sparse

## data\_transformation\_config.py – Defining Paths and Parameters

```
Network_security > entity > config_entity.py > ModelTrainerConfig > __init__
51  class DataTransformationConfig:
52      def __init__(self, training_pipeline_config: TrainingPipelineConfig):
53          self.data_transformation_dir = os.path.join(training_pipeline_config.artifact_dir,
54              training_pipeline.DATA_TRANSFORMATION_DIR_NAME)
55          self.transformed_train_file_path: str = os.path.join(self.data_transformation_dir,
56              training_pipeline.DATA_TRANSFORMATION_TRANSFORMED_DATA_DIR,
57              training_pipeline.TRAIN_FILE_NAME.replace("csv", "npy"))
58          self.transformed_test_file_path: str = os.path.join(self.data_transformation_dir,
59              training_pipeline.DATA_TRANSFORMATION_TRANSFORMED_DATA_DIR,
60              training_pipeline.TEST_FILE_NAME.replace("csv", "npy"))
61          self.transformed_object_file_path: str = os.path.join(self.data_transformation_dir,
62              training_pipeline.DATA_TRANSFORMATION_TRANSFORMED_OBJECT_DIR,
63              training_pipeline.PREPROCESSING_OBJECT_FILE_NAME)
```

self.data\_transformation\_dir → Artifacts/04\_17\_2025\_13\_45\_20/data transformation

self.transformed\_train\_file\_path → Artifacts/04\_17\_2025\_13\_45\_20/data transformation/transformed/train.npy

self.transformed\_test\_file\_path → Artifacts/04\_17\_2025\_13\_45\_20/data transformation/transformed/test.npy

self.transformed\_object\_file\_path → Artifacts/04\_17\_2025\_13\_45\_20/data transformation/transformed\_object/preprocessing.pkl

## utils.py – Saving NumPy Arrays and Python Objects for Pipeline Persistence

```
Network_security > utils > main_utils > utils.py > write_yaml_file
27  def save_numpy_array(file_path: str, array: np.array):
28      try:
29          dir_path = os.path.dirname(file_path)
30          os.makedirs(dir_path, exist_ok=True)
31          with open(file_path, "wb") as file_obj:
32              np.save(file_obj, array)
33      except Exception as e:
34          raise NetworkSecurityException(e, sys)
35
36  def save_object(file_path: str, obj: object) -> None:
37      try:
38          logging.info("Entered the save_object method of MainUtils class")
39          dir_path = os.path.dirname(file_path)
40          os.makedirs(dir_path, exist_ok=True)
41          with open(file_path, "wb") as file_obj:
42              pickle.dump(obj, file_obj)
43      logging.info("Exited the save_object method of MainUtils class")
44      except Exception as e:
45          raise NetworkSecurityException(e, sys)
```

Open the file in binary write mode and saves the NumPy array using np.save(). This is used for .npy file saving

Serializes the object using pickle and writes it to disk in binary format (commonly used for saving models, transformers)

## data\_transformation.py – Generating Transformed Arrays and Preprocessor for Model Training

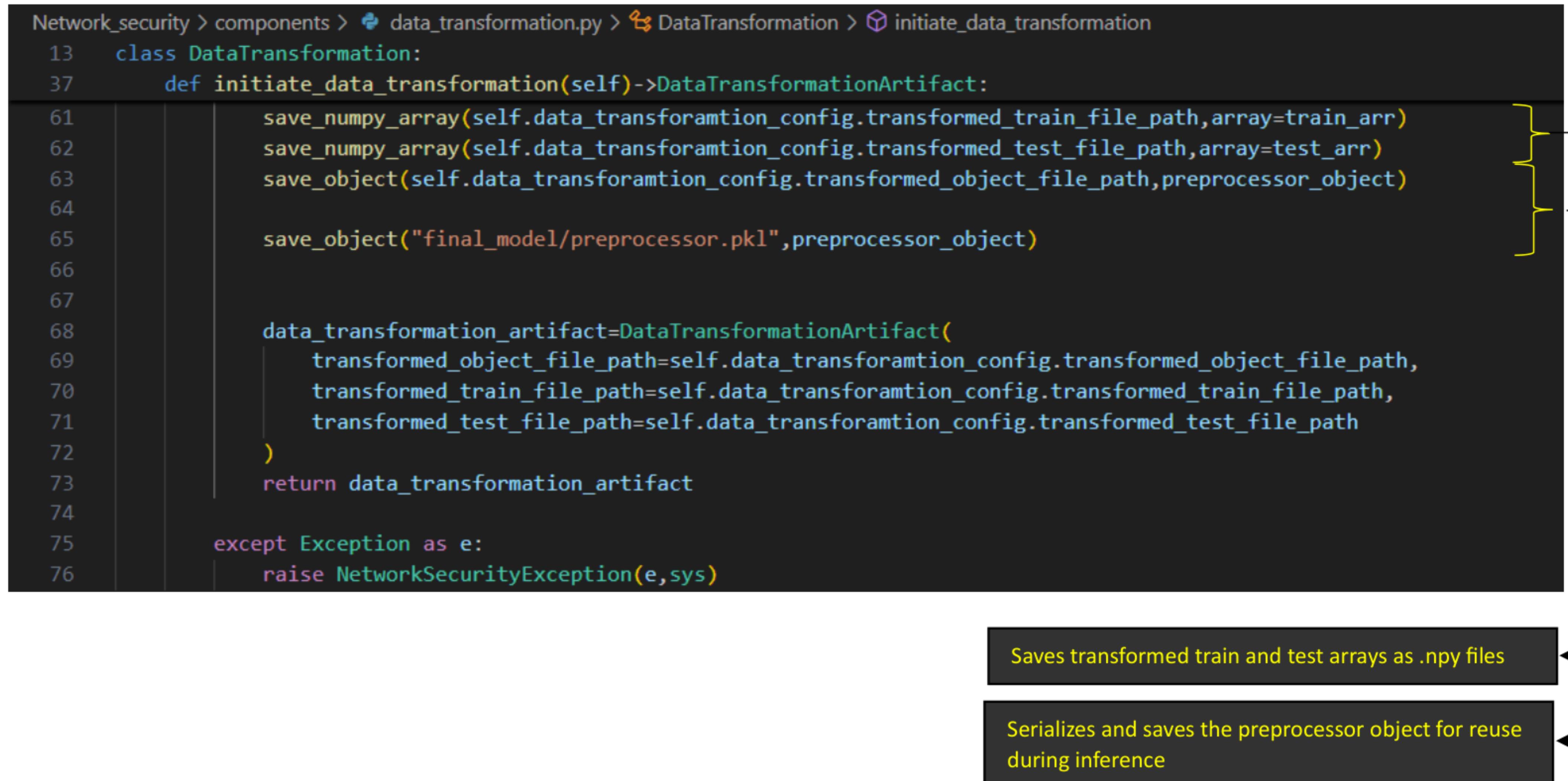
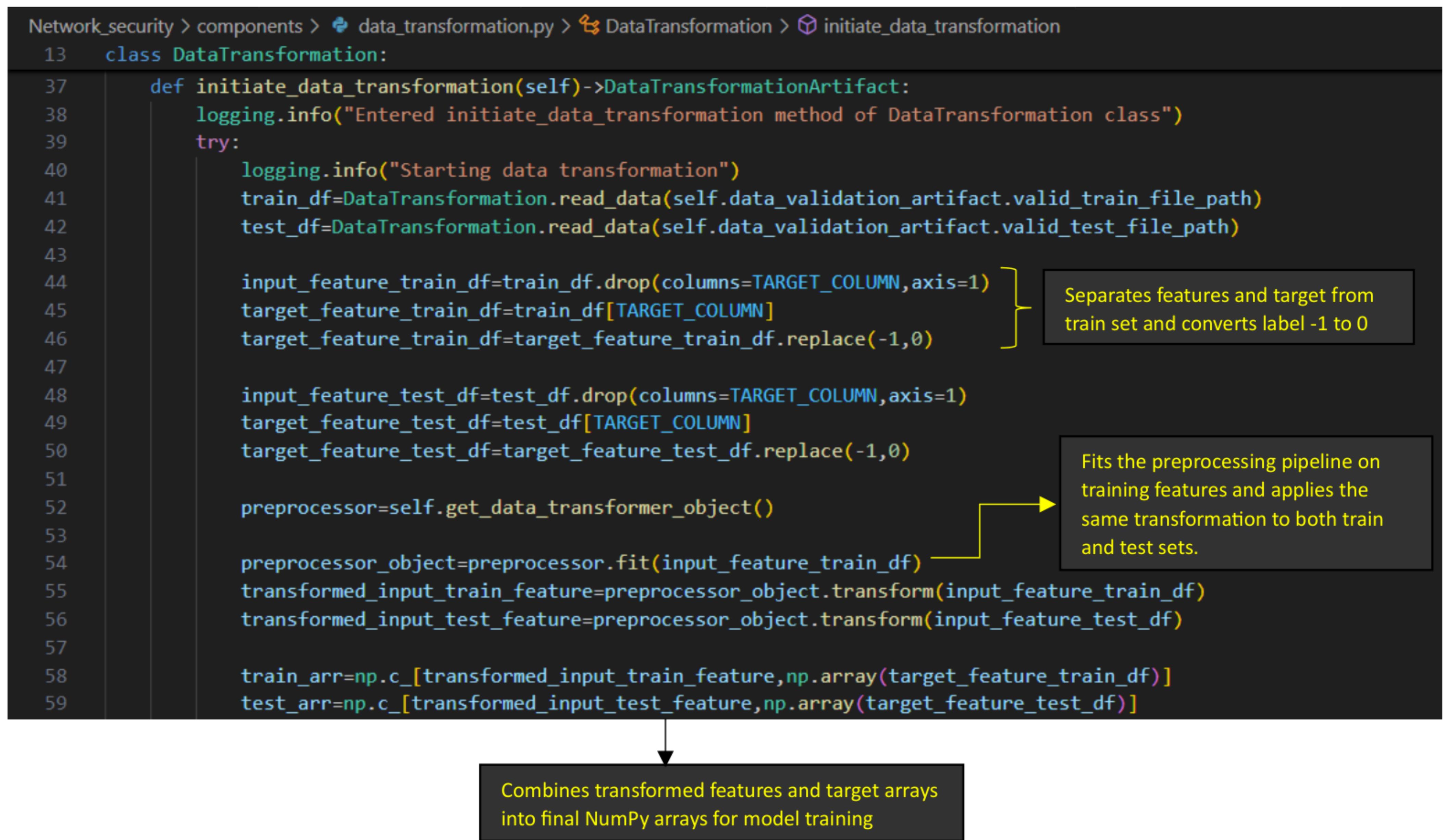
```
Network_security > components > data_transformation.py > DataTransformation > get_data_transformer_object
1 import sys,os
2 import numpy as np
3 import pandas as pd
4 from sklearn.impute import KNNImputer
5 from sklearn.pipeline import Pipeline
6 from Network_security.logging.logger import logging
7 from Network_security.exception.exception import NetworkSecurityException
8 from Network_security.constants.training_pipeline import TARGET_COLUMN
9 from Network_security.constants.training_pipeline import DATA_TRANSFORMATION_IMPUTER_PARAMS
10 from Network_security.entity.artifact_entity import DataTransformationArtifact,DataValidationArtifact
11 from Network_security.entity.config_entity import DataTransformationConfig
12 from Network_security.utils.main_utils.utils import save_numpy_array,save_object
13 class DataTransformation:
14     def __init__(self,data_validation_artifact:DataValidationArtifact,
15                  data_transformation_config:DataTransformationConfig):
16         try:
17             self.data_validation_artifact:DataValidationArtifact=data_validation_artifact
18             self.data_transformation_config:DataTransformationConfig=data_transformation_config
19         except Exception as e:
20             raise NetworkSecurityException(e,sys)
21     @staticmethod
22     def read_data(file_path)->pd.DataFrame:
23         try:
24             return pd.read_csv(file_path)
25         except Exception as e:
26             raise NetworkSecurityException(e,sys)
```

```
Network_security > components > data_transformation.py > DataTransformation > initiate_data_transformation
13 class DataTransformation:
14     def get_data_transformer_object(cls)->Pipeline:
15         logging.info("Entered get_data_transformer_object method of DataTransformation class")
16         try:
17             imputer:KNNImputer=KNNImputer(**DATA_TRANSFORMATION_IMPUTER_PARAMS)
18             logging.info("initialize KNNImputer with {DATA_TRANSFORMATION_IMPUTER_PARAMS}")
19             processor:Pipeline=Pipeline([(("imputer",imputer)])
20             return processor
21         except Exception as e:
22             raise NetworkSecurityException(e,sys)
```

Even though Pipeline currently wraps only KNNImputer, using it ensures safe and consistent preprocessing. It becomes essential when adding steps like StandardScaler or PCA, as it prevents data leakage by fitting only on training data and applying the same logic to test data. This prevents the model from unintentionally learning patterns from the test set during preprocessing.

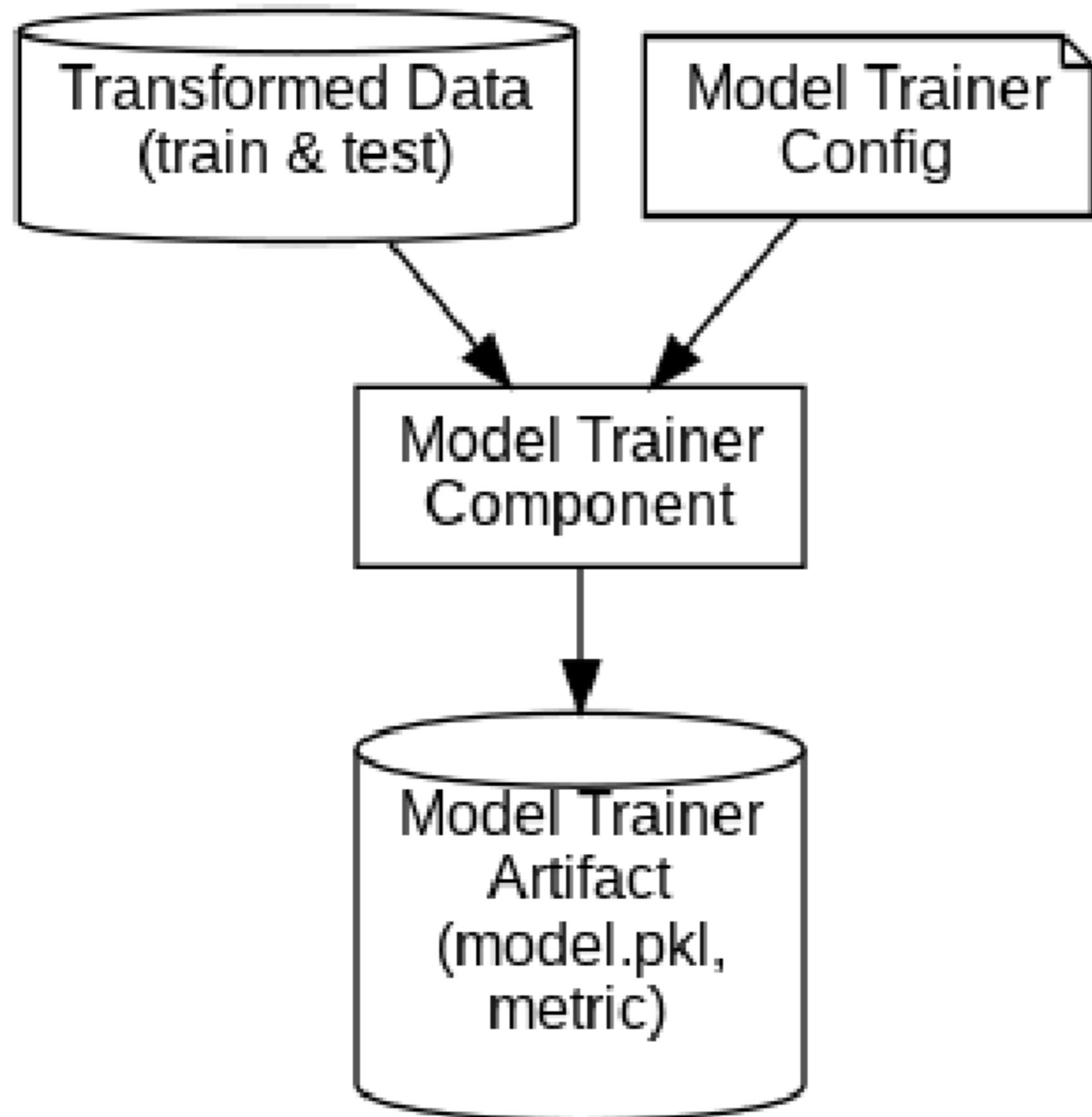
unpacks each key-value pair from the dictionary and passes them as individual keyword arguments to the KNNImputer

## Phishing Website Detection – End-to-End ML Pipeline Project



## Model Trainer Flow

### Training and Saving Model with model.pkl



### training\_pipeline/init.py - Constants for Model Trainer Configuration

```
Network_security > constants > training_pipeline > _init_.py > DATA_TRANSFORMATION_IMPUTER_PARAMS
42 MODEL_TRAINER_DIR_NAME:str="model_trainer"
43 MODEL_TRAINER_TRAINED_MODEL_DIR:str="trained_model"
44 MODEL_TRAINER_TRAINED_MODEL_NAME:str="model.pkl"
45 MODEL_TRAINER_EXPECTED_SCORE:float=0.6
46 MODEL_TRAINER_OVER_FITTING_UNDER_FITTING_THRESOLD:float=0.05
47
48 SAVED_MODEL_DIR=os.path.join("saved_models")
49 MODEL_FILE_NAME="model.pkl"
50
51 TRAINING_BUCKET_NAME="kunalawsbucketsns"
```

### model\_trainer\_config.py – Defining Paths and Parameters for Model Training

```
65 class ModelTrainerConfig:
66     def __init__(self,training_pipleline_config:TrainingPipelineConfig):
67         self.model_trainer_config_dir=os.path.join(
68             training_pipleline_config.artifact_dir,training_pipeline.MODEL_TRAINER_DIR_NAME
69         )
70         self.trained_model_file_path=os.path.join(
71             self.model_trainer_config_dir,training_pipeline.MODEL_TRAINER_TRAINED_MODEL_DIR,
72             training_pipeline.MODEL_TRAINER_TRAINED_MODEL_NAME
73         )
74         self.expected_accuracy:float=training_pipeline.MODEL_TRAINER_EXPECTED_SCORE
75         self.overfitting_underfitting_thresold=training_pipeline.MODEL_TRAINER_OVER_FITTING_UNDER_FITTING_THRESOLD
```

Annotations for the code:

- self.model\_trainer\_config\_dir → Artifacts/04\_17\_2025\_13\_45\_20/model trainer
- self.trained\_model\_file\_path → Artifacts/04\_17\_2025\_13\_45\_20/model trainer/trained\_model/model.pkl

## Phishing Website Detection – End-to-End ML Pipeline Project

### utils.py – Loading Pickle Objects and NumPy Arrays for Reuse in Pipeline

```
Network_security > utils > main_utils > utils.py > load_object
47 def load_object(file_path:str)->object:
48     try:
49         if not os.path.exists(file_path):
50             raise Exception("The file: {file_path} does not exist")
51         with open(file_path,'rb') as file_obj:
52             print(file_obj)
53             return pickle.load(file_obj)
54     except Exception as e:
55         raise NetworkSecurityException(e,sys)
56
57 def load_numpy_array(file_path:str)->np.array:
58     try:
59         with open(file_path,'rb') as file_obj:
60             return np.load(file_obj)
61     except Exception as e:
62         raise NetworkSecurityException(e,sys)
```

```
Network_security > utils > main_utils > utils.py > load_object
64 def evaluate_models(x_train,y_train,x_test,y_test,models,param):
65     try:
66         report={}
67         for i in range(len(list(models))):
68             model=list(models.values())[i]
69             para=param[list(models.keys())[i]]]
70
71             gs=GridSearchCV(model,para,cv=3)
72             gs.fit(x_train,y_train)]]
73
74             model.set_params(**gs.best_params_)
75             model.fit(x_train,y_train)]]
76
77             y_train_pred=model.predict(x_train)
78             y_test_pred=model.predict(x_test)]]
79
80             train_model_score = accuracy_score(y_train, y_train_pred)
81             test_model_score = accuracy_score(y_test, y_test_pred)
82
83             report[list(models.keys())[i]] = test_model_score → Stores the test accuracy of each model in the report dictionary
84
85     return report
86
87     except Exception as e:
88         raise NetworkSecurityException(e,sys)
```

### classification\_report.py – Calculating F1, Precision, and Recall Scores

```
Network_security > utils > ml_utils > metric > classification_report.py > get_classification_score
1 import os,sys
2 from Network_security.entity.artifact_entity import ClassificationMetricArtifact
3 from Network_security.exception.exception import NetworkSecurityException
4 from sklearn.metrics import f1_score,recall_score,precision_score
5
6 def get_classification_score(y_true,y_pred):
7     try:
8         model_f1_score=f1_score(y_true,y_pred)
9         model_recall_score=recall_score(y_true,y_pred)
10        model_precision_score=precision_score(y_true,y_pred)
11
12        classification_metric=ClassificationMetricArtifact(f1_score=model_f1_score,
13                                                               recall_score=model_recall_score,
14                                                               precision_score=model_precision_score)
15
16        return classification_metric
17    except Exception as e:
18        raise NetworkSecurityException(e,sys)]]
```

Calculates F1, recall, and precision scores from predicted and actual labels, then wraps them into a structured artifact for downstream evaluation use

## model\_trainer.py – Training Models with Hyperparameter Tuning and Accuracy Evaluation

```
Network_security > components > model_trainer.py > ModelTrainer > __init__
1  import os,sys
2  from Network_security.exception.exception import NetworkSecurityException
3  from Network_security.logging.logger import logging
4  from Network_security.entity.config_entity import ModelTrainerConfig
5  from Network_security.entity.artifact_entity import DataTransformationArtifact,ModelTrainerArtifact
6  from Network_security.utils.main_utils.utils import save_object,load_object,load_numpy_array,evaluate_models
7  from Network_security.utils.ml_utils.metric.classification_report import get_classification_score
8  from Network_security.utils.ml_utils.model.estimator import NetworkModel
9
10 from sklearn.linear_model import LogisticRegression
11 from sklearn.tree import DecisionTreeClassifier
12 from sklearn.neighbors import KNeighborsClassifier
13 from sklearn.ensemble import RandomForestClassifier,AdaBoostClassifier,GradientBoostingClassifier
14 import mlflow

18 class ModelTrainer:
19     def __init__(self,model_trainer_config:ModelTrainerConfig,data_transforamtion_artifact:DataTransformationArtifact):
20         try:
21             self.model_trainer_config=model_trainer_config
22             self.data_transforamtion_artifact=data_transforamtion_artifact
23         except Exception as e:
24             raise NetworkSecurityException(e,sys)
25
26     def track_mlflow(self,best_model,classification_train_metric):
27         with mlflow.start_run():
28             f1_score=classification_train_metric.f1_score
29             precision_score=classification_train_metric.precision_score
30             recall_score=classification_train_metric.recall_score
31
32             mlflow.log_metric("f1_score",f1_score)
33             mlflow.log_metric("precision_score",precision_score)
34             mlflow.log_metric("recall_score",recall_score)
35             mlflow.sklearn.log_model(best_model,"model")
```

```
Network_security > components > model_trainer.py > ModelTrainer > track_mlflow
18 class ModelTrainer:
19     def train_model(self,x_train,y_train,x_test,y_test):
20         models={
21             "Random Forest":RandomForestClassifier(verbose=1),
22             "Decision Tree":DecisionTreeClassifier(),
23             "Gradient Boosting":GradientBoostingClassifier(verbose=1),
24             "Logistic Regression":LogisticRegression(verbose=1),
25             "AdaBoost":AdaBoostClassifier()
26         }
27         params={
28             "Decision Tree": {
29                 'criterion':['gini', 'entropy', 'log_loss'],
30                 # 'splitter':['best','random'],
31                 # 'max_features':['sqrt','log2'],
32             },
33             "Random Forest":{
34                 # 'criterion':['gini', 'entropy', 'log_loss'],
35
36                 # 'max_features':['sqrt','log2',None],
37                 'n_estimators': [8,16,32,128,256]
38             },
39             "Gradient Boosting":{
40                 # 'loss':['log_loss', 'exponential'],
41                 'learning_rate':[.1,.01,.05,.001],
42                 'subsample':[0.6,0.7,0.75,0.85,0.9],
43
44                 # 'criterion':['squared error', 'friedman mse'],
45             }
46         }
```

Phishing Website Detection – End-to-End ML Pipeline Project

Network\_security > components > model\_trainer.py > ModelTrainer > train\_model

```
18  class ModelTrainer:  
37      def train_model(self,x_train,y_train,x_test,y_test):  
71          model_report:dict=evaluate_models(x_train=x_train,y_train=y_train,x_test=x_test,y_test=y_test,  
72                                         models=models,param=params) →  
73          best_model_score=max(sorted(model_report.values()))  
74          best_model_name=list(model_report.keys())[list(model_report.values()).index(best_model_score)] ↓  
75          best_model=models[best_model_name]  
76          y_train_pred=best_model.predict(x_train)  
77          classification_train_metric=get_classification_score(y_true=y_train,y_pred=y_train_pred)  
78          # Track the mlflow  
81          self.track_mlflow(best_model,classification_train_metric) ← Logs training metrics and model into MLflow for tracking and versioning  
82          y_test_pred=best_model.predict(x_test)  
83          classification_test_metric=get_classification_score(y_true=y_test,y_pred=y_test_pred)  
84          self.track_mlflow(best_model,classification_test_metric)  
85          preprocessor=load_object(file_path=self.data_transformation_artifact.transformed_object_file_path) ↓  
86          model_dir_path=os.path.dirname(self.model_trainer_config.trained_model_file_path) ↓  
87          os.makedirs(model_dir_path,exist_ok=True)  
88          Network_model=NetworkModel(preprocessor=preprocessor,model=best_model)  
89          save_object(self.model_trainer_config.trained_model_file_path,obj=Network_model) ↓  
90          save_object("final_model/model.pkl",best_model) ↓  
91          ↓  
92          Wraps both preprocessor and model into a custom class, then saves it. Also  
93          saves the plain model separately for reference
```

Network\_security > components > model\_trainer.py > ModelTrainer > train\_model

```
18  class ModelTrainer:  
37      def train_model(self,x_train,y_train,x_test,y_test):  
93          model_trainer_artifact=ModelTrainerArtifact(trained_model_file_path=  
94                                         self.model_trainer_config.trained_model_file_path,  
95                                         train_metric_artifact=classification_train_metric,  
96                                         test_metric_artifact=classification_test_metric)  
97  
98          logging.info(f"Model trainer artifact: {model_trainer_artifact}")  
99          return model_trainer_artifact  
100  
101     def initiate_model_trainer(self)->ModelTrainerArtifact:  
102         try:  
103             trainer_file_path=self.data_transformation_artifact.transformed_train_file_path  
104             test_file_path=self.data_transformation_artifact.transformed_test_file_path  
105  
106             train_arr=load_numpy_array(trainer_file_path)  
107             test_arr=load_numpy_array(test_file_path) } ↓  
108             x_train,y_train,x_test,y_test=(  
109                 train_arr[:, :-1],  
110                 train_arr[:, -1],  
111                 test_arr[:, :-1],  
112                 test_arr[:, -1]) } ↓  
113             ↓  
114             ↓  
115             Loads the NumPy arrays containing transformed  
116             features and target from disk  
117             ↓  
118             Splits the arrays into input features (x) and  
119             target labels (y) for both train and test sets
```

## **artifact\_entity.py –Artifacts for Transformation, Metrics, and Model Trainer Output**

```
23 @dataclass
24 class ClassificationMetricArtifact:
25     f1_score:float
26     precision_score:float
27     recall_score:float
28
29 @dataclass
30 class ModelTrainerArtifact:
31     trained_model_file_path:str
32     train_metric_artifact:ClassificationMetricArtifact
33     test_metric_artifact:ClassificationMetricArtifact
```

## Create an IAM User

### Foundation for Secure and Automated AWS Access

The very first step when working with AWS programmatically is to create an **IAM (Identity and Access Management) user**. This user allows secure, credential-based access to AWS services through CLI tools, SDKs like boto3, or automation pipelines.

### Creating an IAM User

- Navigate to the AWS Management Console → IAM → Users → Add user
- Enable Programmatic access to generate the Access Key ID and Secret Access Key
- Assign required permissions: For full access: attach the AdministratorAccess policy

The screenshot shows the AWS IAM User Details page for a user named 'testsecurity'. The 'Summary' tab is selected, displaying the ARN, creation date (April 14, 2025), and access status (Console access Disabled). The 'Permissions' tab is also visible, showing the attached 'Permissions policies (1)' which is the 'AdministratorAccess' policy. The policy details show it is an AWS managed - job function policy attached directly.

## Why AWS S3?

### Making ML Pipeline Cloud-Ready and Production-Friendly

AWS S3 is used in this project to store important pipeline artifacts like the trained model (model.pkl), preprocessing object (preprocessor.pkl), and drift reports securely in the cloud.

This solves three key problems:

1. **Storage** – Keeps your model and artifacts safe even if your local environment is cleaned or crashes.
2. **Access** – Makes your trained models accessible from anywhere (e.g., EC2, Lambda, or other APIs).
3. **Deployment-Ready** – Cloud storage like S3 is essential when moving towards real-time deployments, CI/CD pipelines, or serverless architecture.

By syncing both the artifact directory and final model to S3, the pipeline becomes portable, scalable, and ready for production.

## s3\_syncer.py – Syncing Artifacts like `preprocessor.pkl` to AWS S3

This file marks the beginning of the cloud integration step.

The S3Sync class defines a method sync\_folder\_to\_s3, which uses the AWS CLI to upload folders (like final models, drift reports, and preprocessing objects) to a specified S3 bucket for storage, sharing, or deployment.

```
Network_security > cloud > s3_syncer.py > ...
1 import os
2
3 class S3Sync:
4     def sync_folder_to_s3(self, folder, aws_bucket_url):
5         command=f"aws s3 sync {folder} {aws_bucket_url}"
6         os.system(command)
```

Defines a method that takes the local folder path and target S3 bucket URL as inputs

Builds the AWS CLI command to sync the local folder with the specified S3 bucket

## main.py – Running the Complete Training Pipeline and Syncing Artifacts to S3

```
Network_security > pipeline > training_pipeline.py > ...
1 import sys,os
2 from Network_security.exception.exception import NetworkSecurityException
3 from Network_security.logging.logger import logging
4 from Network_security.components.data_ingestion import DataIngestion
5 from Network_security.components.data_validation import DataValidation
6 from Network_security.components.data_transformation import DataTransformation
7 from Network_security.components.model_trainer import ModelTrainer
8 from Network_security.entity.config_entity import(
9     TrainingPipelineConfig,DataIngestionConfig,DataValidationConfig,DataTransformationConfig,ModelTrainerConfig
10 )
11 from Network_security.entity.artifact_entity import(
12     DataIngestionArtifact,DataValidationArtifact,DataTransformationArtifact,ModelTrainerArtifact
13 )
14
15 from Network_security.constants.training_pipeline import TRAINING_BUCKET_NAME
16 from Network_security.cloud.s3_syncer import S3Sync
17 from Network_security.constants.training_pipeline import SAVED_MODEL_DIR
18
19 class TrainingPipeline:
20     def __init__(self):
21         self.training_pipeline_config=TrainingPipelineConfig()
22         self.s3_sync = S3Sync()
```

Creates an instance of the S3Sync class to handle syncing local folders to AWS S3

```
Network_security > pipeline > training_pipeline.py > ...
19 class TrainingPipeline:
20
21     def start_data_ingestion(self):
22         try:
23             self.data_ingestion_config=DataIngestionConfig(training_pipeline_config=self.training_pipeline_config)
24             logging.info("Initiate data ingestion")
25             data_ingestion=DataIngestion(data_ingestion_config=self.data_ingestion_config)
26             data_ingestion_artifact=data_ingestion.initiate_data_ingestion()
27             logging.info(f"Data ingestion completed and artifact:{data_ingestion_artifact}")
28             return data_ingestion_artifact
29         except Exception as e:
30             raise NetworkSecurityException(e,sys)
31
32     def start_data_validation(self,data_ingestion_artifact:DataIngestionArtifact):
33         try:
34             data_validation_config>DataValidationConfig(training_pipeline_config=self.training_pipeline_config)
35             logging.info("Initiate data validation")
36             data_validation=DataValidation(data_ingestion_artifact=data_ingestion_artifact,
37                                         data_validation_config=data_validation_config)
38             dat_validation_artifact=data_validation.initiate_data_validation()
39             logging.info(f"Data validation completed and artifact:{dat_validation_artifact}")
40             return dat_validation_artifact
41         except Exception as e:
42             raise NetworkSecurityException(e,sys)
```

Runs the ingestion process and returns a DataIngestionArtifact containing paths to the train and test CSV files

Validates the ingested data and returns a DataValidationArtifact with valid/invalid file paths and drift report

## Phishing Website Detection – End-to-End ML Pipeline Project

```
Network_security > pipeline > 📁 training_pipeline.py > ...
19   class TrainingPipeline:
47     def start_data_transformation(self,data_validation_artifact:DataValidationArtifact):
48       try:
49         data_transformation_config=DataTransformationConfig(training_pipeline_config=self.training_pipeline_config)
50         logging.info("Initiate data transformation")
51         data_transformation>DataTransformation(data_transformation_config=data_transformation_config,
52                                               data_validation_artifact=data_validation_artifact)
53         data_transformation_artifact=data_transformation.initiate_data_transformation()
54         logging.info(f"Data transformation completed and artifact:{data_transformation_artifact}")
55         return data_transformation_artifact
56     except Exception as e:
57       raise NetworkSecurityException(e,sys)   Transforms the data and returns a DataTransformationArtifact
58                                         with transformed arrays and the saved preprocessor
59
60     def start_model_trainer(self,data_transformation_artifact:DataTransformationArtifact):
61       try:
62         model_trainer_config=ModelTrainerConfig(training_pipeline_config=self.training_pipeline_config)
63         logging.info("Initiate model trainer")
64         model_trainer=ModelTrainer(model_trainer_config=model_trainer_config,
65                                   data_transformation_artifact=data_transformation_artifact)
66         model_trainer_artifact=model_trainer.initiate_model_trainer()
67         logging.info(f"Model trainer completed and artifact:{model_trainer_artifact}")
68         return model_trainer_artifact
69     except Exception as e:
69       raise NetworkSecurityException(e,sys)   Trains and evaluates models, returning a ModelTrainerArtifact
                                         with model path and training/testing metrics
```

```
Network_security > pipeline > 📁 training_pipeline.py > 🏃 TrainingPipeline > ⚙️ run_pipeline
19   class TrainingPipeline:
73     def sync_artifact_dir_to_s3(self):
74       try:
75         aws_bucket_url = f"s3://{TRAINING_BUCKET_NAME}/artifact/{self.training_pipeline_config.timestamp}"
76         self.s3_sync.sync_folder_to_s3(folder = self.training_pipeline_config.artifact_dir,aws_bucket_url=aws_bucket_url)
77     except Exception as e:
78       raise NetworkSecurityException(e,sys)   Builds the target S3 path for syncing all generated artifacts using the pipeline
79
80     def sync_saved_model_dir_to_s3(self):
81       try:
82         aws_bucket_url = f"s3://{TRAINING_BUCKET_NAME}/final_model/{self.training_pipeline_config.timestamp}"
83         self.s3_sync.sync_folder_to_s3(folder = self.training_pipeline_config.model_dir,aws_bucket_url=aws_bucket_url)
84     except Exception as e:
85       raise NetworkSecurityException(e,sys)   Uploads the entire local artifact directory to the specified S3 path, and syncs the
86
86     def run_pipeline(self):
87       try:
88         print(f"■ Local artifact folder: {self.training_pipeline_config.artifact_dir}")
89         print(f"■ Local model folder: {self.training_pipeline_config.model_dir}")
90         data_ingestion_artifact=self.start_data_ingestion()
91         data_validation_artifact=self.start_data_validation(data_ingestion_artifact=data_ingestion_artifact)
92         data_transformation_artifact=self.start_data_transformation(data_validation_artifact=data_validation_artifact)
93         model_trainer_artifact=self.start_model_trainer(data_transformation_artifact=data_transformation_artifact)
94         self.sync_artifact_dir_to_s3()
95         self.sync_saved_model_dir_to_s3()
96         return model_trainer_artifact
97     except Exception as e:
97       raise NetworkSecurityException(e,sys)
```

## S3 Bucket Overview – Uploaded Artifacts and Final Model from Training Pipeline

kunalawsbucketsns [Info](#)

Name	Type	Last modified	Size	Storage class
<a href="#">artifact/</a>	Folder	-	-	-
<a href="#">final_model/</a>	Folder	-	-	-

This screenshot shows the contents of your S3 bucket named **kunalawsbucketsns**, where two folders—`artifact/` and `final_model/`—have been successfully uploaded. These folders are automatically created and synced to AWS S3 at the end of the training pipeline. The `artifact/` folder contains all intermediate outputs like ingested data, validation reports, transformed files, and other generated artifacts during each pipeline stage. The `final_model/` folder contains the serialized trained model and preprocessing object, making them easily accessible for deployment. This sync happens as part of the `run_pipeline()` method, which triggers two S3 upload operations through the methods `sync_artifact_dir_to_s3()` and `sync_saved_model_dir_to_s3()` using the AWS CLI in the background.

## table.html – Rendering Prediction Results in a Browser-Friendly Table

```
templates > table.html > html
1  <!DOCTYPE html> } Declares the document as an HTML5 webpage
2  <html>
3  <head>
4      <title>Data Table</title> } Sets the browser tab title to "Data Table"
5      <style>
6          table {
7              border-collapse: collapse;
8              width: 100%;
9          }
10         th,td {
11             border: 1px solid black;
12             padding: 8px;
13             text-align: left;
14         }
15     </style>
16 </head>
17 <body>
18     <h2>Predicted Data</h2> } Shows a heading above the table
19     {{table | safe}} → This is a Jinja2 placeholder. FastAPI replaces this with the actual HTML table
20     generated from the pandas DataFrame in your backend.
21     safe part tells Jinja2: "Don't escape this as text, render it as real HTML."
22 </body>
23 </html>
```

## app.py – FastAPI Interface for Training and Prediction

This file creates a FastAPI application with routes to trigger the training pipeline and generate predictions. It loads the trained model and preprocessor, accepts a CSV input, and returns predictions in a browser-rendered HTML table. The app integrates Jinja2 for templating and supports cross-origin requests for deployment.

```
app.py > ...
1 import os,sys
2 import certifi
3 ca=certifi.where()
4 from dotenv import load_dotenv
5 load_dotenv()
6 mongo_db_url=os.getenv("MONGODB_URL_KEY")
7 print(mongo_db_url)
8 import pymongo
9 from Network_security.exception.exception import NetworkSecurityException
10 from Network_security.logging.logger import logging
11 from Network_security.pipeline.training_pipeline import TrainingPipeline
12 from fastapi.middleware.cors import CORSMiddleware
13 from fastapi import FastAPI,Request,UploadFile,File
14 from uvicorn import run as app_run
15 from fastapi.responses import Response
16 from starlette.responses import RedirectResponse
17 import pandas as pd
18
19 from Network_security.utils.main_utils.utils import load_object
20 from Network_security.utils.ml_utils.model.estimator import NetworkModel
21
22
23 from Network_security.constants.training_pipeline import DATA_INGESTION_DATABASE_NAME,DATA_INGESTION_COLLECTION_NAME
```

```
app.py > ...
25 client=pymongo.MongoClient(mongo_db_url,tlsCAFile=ca) → Establishes a secure MongoDB connection using certifi for TLS validation
26 database=client[DATA_INGESTION_DATABASE_NAME]
27 collection=database[DATA_INGESTION_COLLECTION_NAME]
28
29 app=FastAPI() → Creates the FastAPI app instance
30 origins=["*"] → Allows any frontend (from any domain) to access your API. The * means "allow all"
31 app.add_middleware(
32     CORSMiddleware,
33     allow_origins=origins,
34     allow_credentials=True,
35     allow_methods=["*"],
36     allow_headers=["*"]
37 ) → Adds CORS middleware (Cross-Origin Resource Sharing).
          This is required when your frontend (like React, Streamlit, or Postman) calls your
          backend API from another domain.
          Without this, browsers block external requests for security reasons.
38
39 from fastapi.templating import Jinja2Templates
40 templates = Jinja2Templates(directory="../templates") → Tells FastAPI to look for HTML files (like
          table.html) inside the templates/ folder
41
42 @app.get("/",tags=["authentication"])
43 async def index():
44     return RedirectResponse(url="/docs") → / redirects to the /docs page, which is the Swagger UI auto-
          generated by FastAPI — used for API testing
```

Jinja2 is a powerful templating engine used to generate HTML dynamically by inserting data from the backend into predefined template files. It is commonly used with web frameworks like Flask and FastAPI

## Phishing Website Detection – End-to-End ML Pipeline Project

```
app.py > index
46 @app.get("/train") → Initializes the entire pipeline class
47 async def train_route():
48     try:
49         train_pipeline=TrainingPipeline()
50         train_pipeline.run_pipeline()
51         return Response("Training is successful")
52     except Exception as e:
53         raise NetworkSecurityException(e,sys)
54
55 @app.post("/predict")
56 async def predict_route(request: Request,file: UploadFile = File(...)): → This route handles POST requests (with file upload).
57     try:
58         df=pd.read_csv(file.file)
59         preprocesor=load_object("final_model/preprocessor.pkl")
60         final_model=load_object("final_model/model.pkl")
61         network_model = NetworkModel(preprocessor=preprocesor,model=final_model)
62         print(df.iloc[0])
63         y_pred = network_model.predict(df)
64         print(y_pred)
65         df['predicted_column'] = y_pred
66         print(df['predicted_column'])
67         df.to_csv('prediction_output/output.csv') → Converts the DataFrame into HTML
68         table_html = df.to_html(classes='table table-striped') → using Bootstrap-like table styling
69         return templates.TemplateResponse("table.html", {"request": request, "table": table_html})
70     except Exception as e:
71         raise NetworkSecurityException(e,sys)
```

```
app.py > predict_route
73 if __name__=="__main__":
74     app_run(app,host="0.0.0.0",port=8080) } Runs your FastAPI app using Uvicorn on port 8080.
75
76 print(" Deployment trigger") 0.0.0.0 means it's accessible externally (e.g., from browser, Postman, or another server)
```

## Dockerfile

### Containerizing the FastAPI Application for AWS Deployment

This Dockerfile defines how the entire machine learning application — including the FastAPI app, model, and all dependencies — is packaged into a lightweight container image, ready for deployment on any platform such as AWS EC2

```
Dockerfile > ...
1 FROM python:3.10-slim-buster → Uses a lightweight official Python 3.10 base image to keep the container small and efficient
2 WORKDIR /app } Sets the working directory to /app and copies all local project files into that directory inside the container
3 COPY . /app
4
5 RUN apt update -y && apt install awscli -y → Updates the package manager and installs the AWS CLI inside the
6                                         container, enabling interaction with AWS services like S3
7 RUN apt-get update && pip install -r requirements.txt → Installs all Python dependencies listed in requirements.txt
8
9 CMD ["python3", "app.py"]
```

↓

Defines the default command that will run when the container starts.  
This launches the FastAPI app (app.py) using Python

## GitHub

### Pushing Local Machine Learning Project to GitHub

```
cd path/to/your/project  
git init  
git add .  
git commit -m "Commit message"  
git remote add origin https://github.com/KunalPShukla/Network_Security.git  
git branch -M main  
git push -u origin main
```

It is good practice to regularly push changes to GitHub to maintain version control and project backups

Files and folders listed in .gitignore will not be pushed to the repository

The screenshot shows a GitHub repository page for 'Network\_Security'. At the top, there's a summary bar with 'main' (branch), '1 Branch', '0 Tags', a search bar ('Go to file'), an 'Add file' button, and a 'Code' button. Below this is a list of commits:

Author	Message	Date
KunalPShukla	Triggering CI/CD pipeline ✓	dc59f33 · last week 19 Commits
	.github/workflows	Final deployment
	Network_security	Cloud s3 data storage
	Network_security_data	Project Structure
	__pycache__	Cloud s3 data storage
	data_schema	The first commit
	final_model	Cloud s3 data storage
	mlruns/0	Cloud s3 data storage
	prediction_output	Pipeline is completed.
	templates	Pipeline is completed.
	valid_data	Pipeline is completed.
	.gitignore	Project Structure

## GitHub Secrets – Securely Storing AWS Credentials for CI/CD Pipeline

It ensures that all sensitive AWS credentials (like access keys and login URIs) are safely stored and can be accessed inside the workflow without exposing them in the code

The screenshot shows the GitHub Secrets interface. On the left, there's a sidebar with options like Branches, Tags, Rules, Actions, Webhooks, Environments, Codespaces, and Pages. Below that is a Security section with Code Security, Deploy keys, and Secrets and variables (which is expanded). Under Actions, Codespaces, and Dependabot are listed. Integrations include GitHub Apps and Email notifications. The main area has tabs for Secrets and Variables, with Secrets selected. It shows two sections: Environment secrets (empty) and Repository secrets. The Repository secrets table lists five secrets:

Name	Last updated	Actions
AWS_ACCESS_KEY_ID	last week	
AWS_ECR_LOGIN_URI	last week	
AWS_REGION	last week	
AWS_SECRET_ACCESS_KEY	last week	
ECR_REPOSITORY_NAME	last week	

A green "New repository secret" button is located at the top right of the Repository secrets section.

## Amazon ECR

### Private Repository for ML App Docker Image

This private ECR repository named networksecurity stores the Docker image of the ML application. The image is built via GitHub Actions and pushed to this repository, from where it is later pulled by an EC2 instance for deployment.

The screenshot shows the Amazon ECR console under the Private registry section. The sidebar includes options for Private registry (Repositories, Features & Settings) and Public registry (Repositories, Settings). The main area shows "Private repositories (1)". A search bar says "Search by repository substring". A table lists the repository "networksecurity" with details: Repository name (networksecurity), URI (redacted), Created at (April 15, 2025, 22:47:43 (UTC+05.5)), Tag immutability (Mutable), and Encryption type (AES-256). Action buttons include View push commands, Delete, Actions, and Create repository.

## GitHub Actions Workflow

# CI/CD Pipeline for Dockerized ML App to AWS ECR & EC2

This GitHub Actions workflow automates the **build, test, and deployment** process of your machine learning FastAPI app

### Continuous Integration (CI)

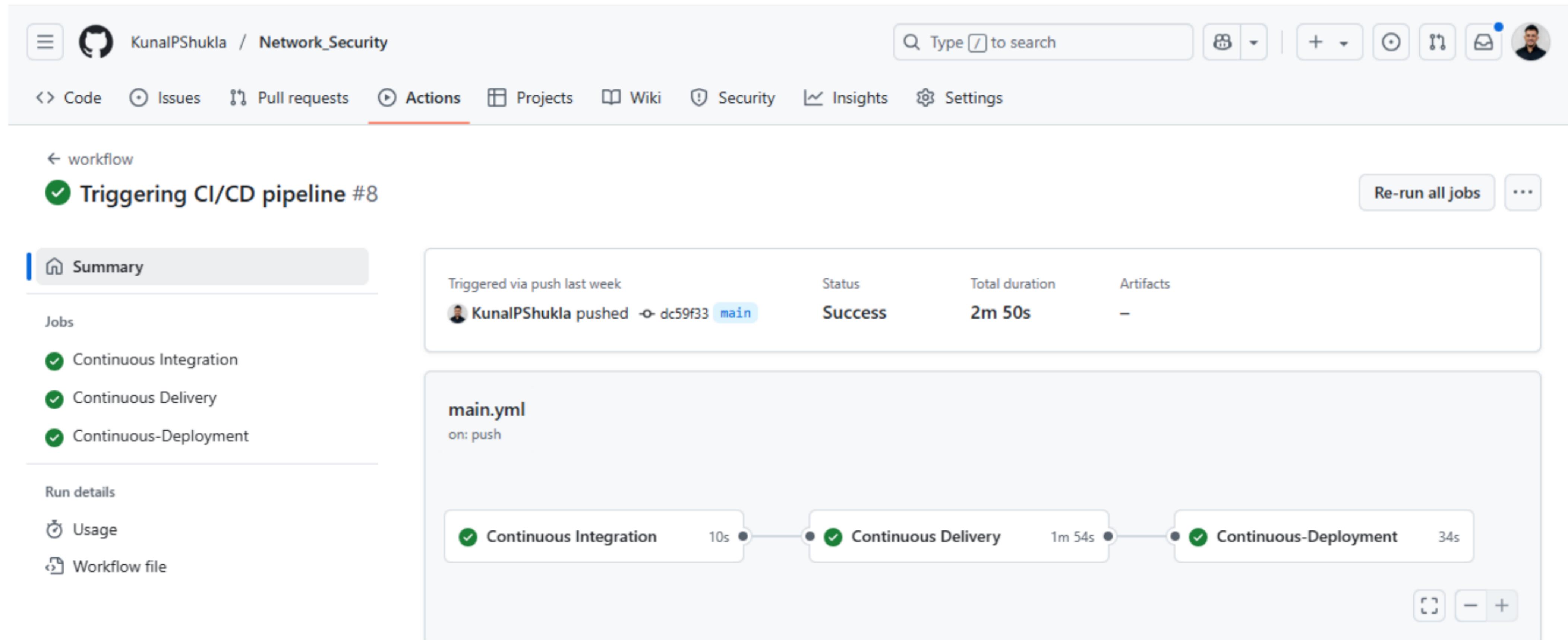
- Linting
- Unit testing
- Runs when a push is made to the main branch

### Continuous Delivery (CD)

- Builds and tags a Docker image
- Pushes it to Amazon ECR

### Continuous Deployment (CD)

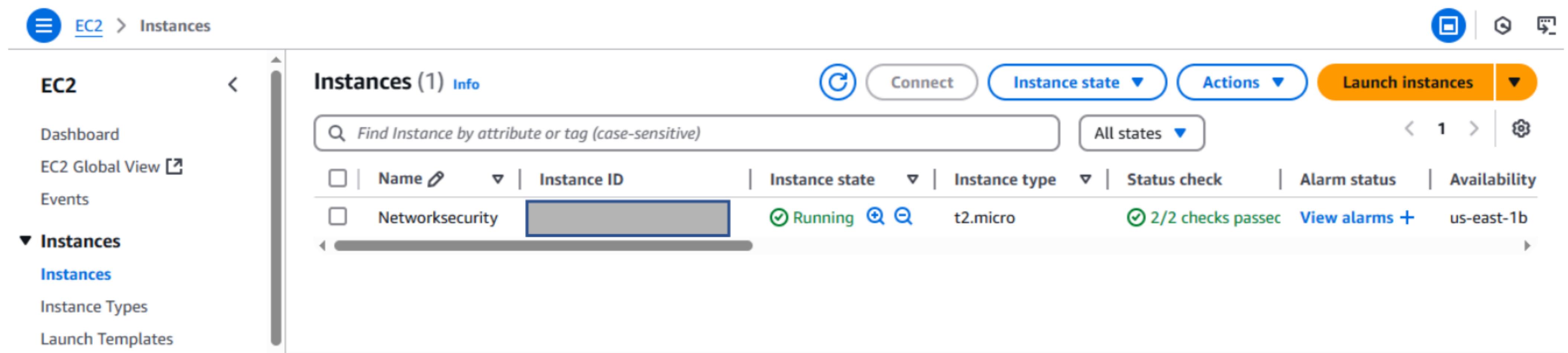
- Runs on a self-hosted EC2 runner
- Pulls the latest image from ECR
- Removes existing containers if running
- Deploys the new image
- Cleans up old images and containers



## Amazon EC2

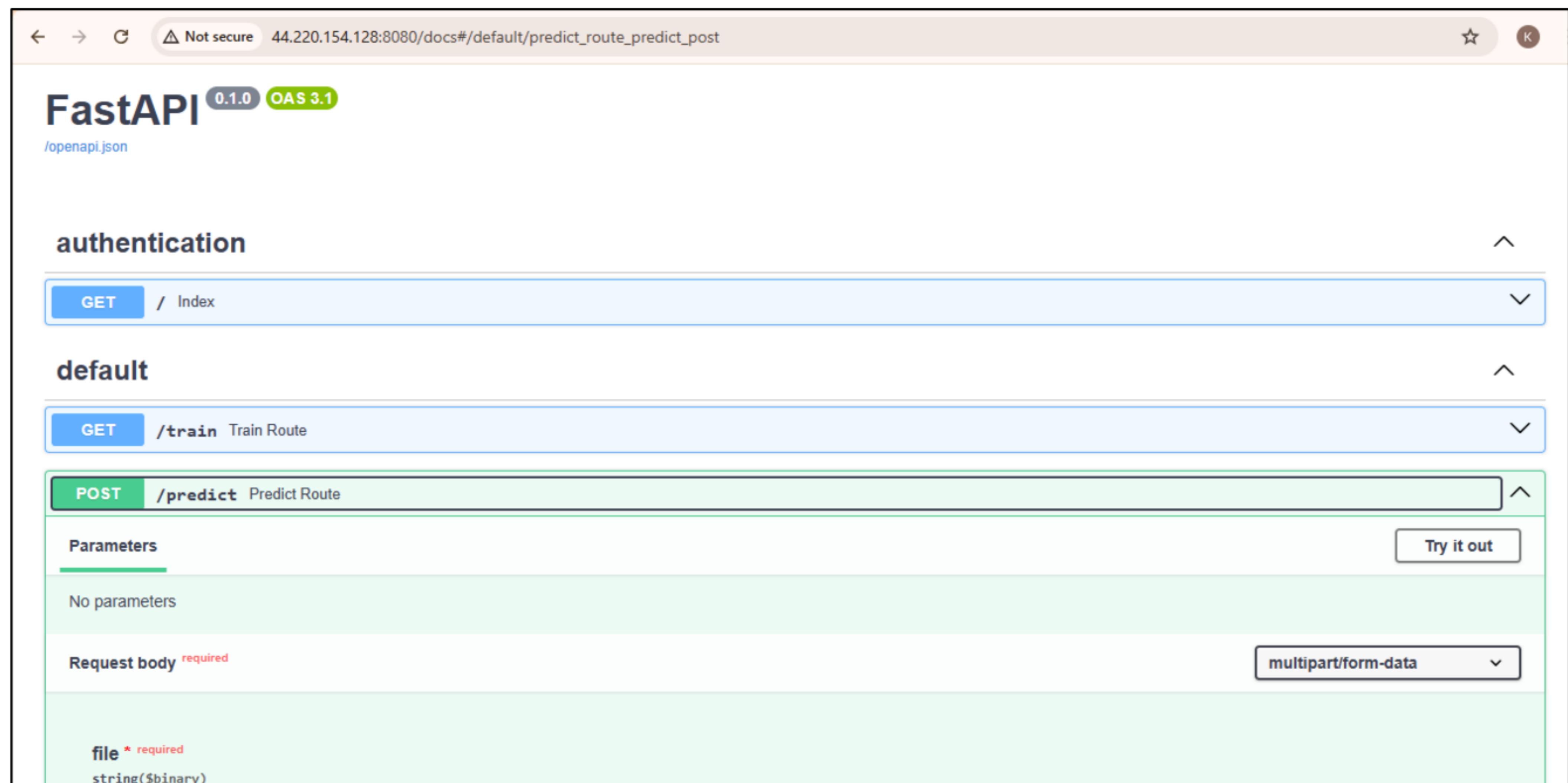
### Hosting the Deployed ML Application

This EC2 instance acts as the **host machine** for running the Dockerized ML application. After the Docker image is pushed to **Amazon ECR** through GitHub Actions, this instance **pulls the image** and starts the containerized app — making it accessible via the public IP



The screenshot shows the AWS EC2 Instances page. On the left, there's a sidebar with 'EC2' selected. The main area has a heading 'Instances (1) Info'. Below it is a search bar and a table with columns: Name, Instance ID, Instance state, Instance type, Status check, Alarm status, and Availability. One row is visible, showing 'Networksecurity' as the name, 'Running' as the state, 't2.micro' as the type, '2/2 checks passed' as the status, and 'us-east-1b' as the availability zone.

This marks the successful deployment of the entire ML pipeline into a **cloud-accessible production environment**.



The screenshot shows the FastAPI documentation interface. At the top, it says 'FastAPI 0.1.0 OAS 3.1'. Below that, there are sections for 'authentication' and 'default'. Under 'default', there are two endpoints: a 'GET /train' endpoint for training routes, and a 'POST /predict' endpoint for predicting routes. The 'POST /predict' endpoint is expanded, showing 'Parameters' (none listed), 'Request body' (required), and a file input field labeled 'file \* required string(\$binary)'. There's also a 'Try it out' button.

### Credits:

This end-to-end machine learning deployment project structure, workflow, and guidance is inspired by **Krish Naik's Udemy course** on scalable ML pipelines and cloud deployment. Special thanks for the foundational knowledge and project framework.