

Arjun Anand Mallya 23110039

Venkatakrishnan E 23110357

Importing the Breast cancer Dataset

```
main.ipynb × workflow.py × serving.py × secret.txt × load_data.py × trainer.py ×
1 import mlrun
2 from sklearn.datasets import load_breast_cancer
3 import pandas as pd
4
5 @mlrun.handler(outputs=["dataset", "label_column"])
6 def breast_cancer_loader(context, format="csv"):
7
8     data = load_breast_cancer(as_frame=True)
9     breast_cancer_df = data.frame
10    breast_cancer_df['target'] = data.target
11
12    context.logger.info('Saving breast cancer dataset to {}'.format(context.artifact_path))
13    context.log_dataset('breast_cancer_dataset', df=breast_cancer_df, format=format, index=False)
14
15    return breast_cancer_df, "target"
16
17 if __name__ == "__main__":
18     with mlrun.get_or_create_ctx("breast_cancer_loader", upload_artifacts=True) as context:
19         breast_cancer_loader(context, context.get_param("format", "csv"))
20
```

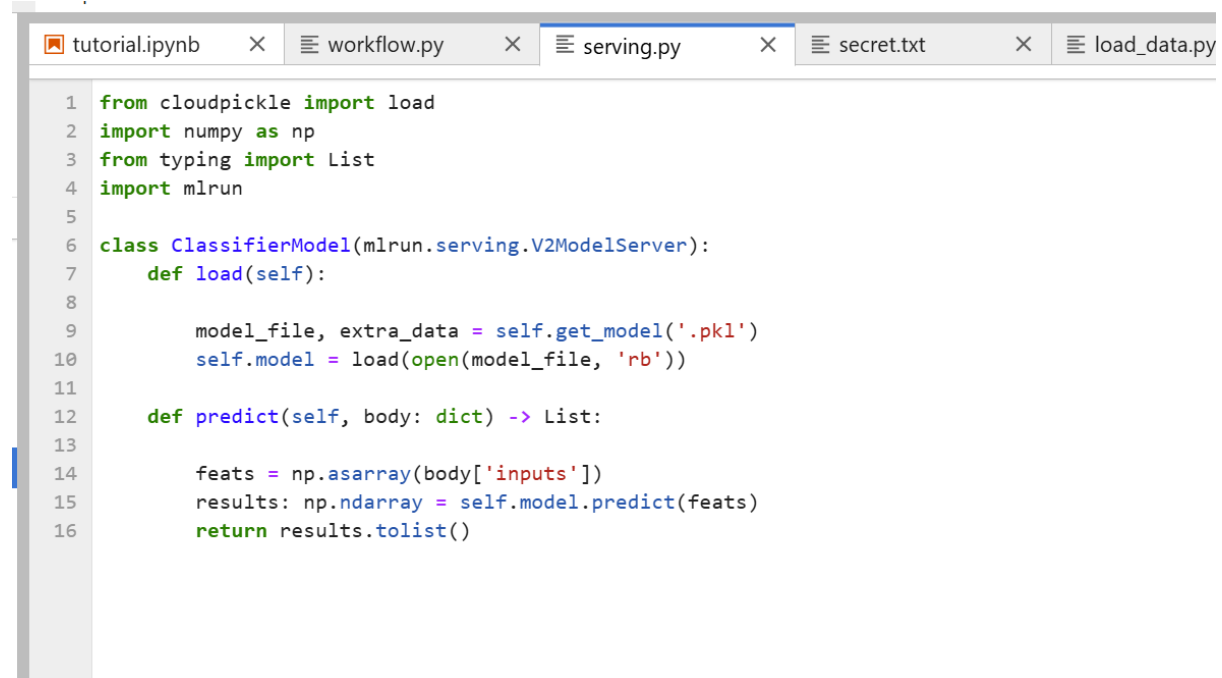
Splitting the data into train/test (Trainer.py)

The dataset was split into 90% train and 10% split. The dataset was loaded using ML Run and prepared for training. It trains a Random Forest classifier.

```
main.ipynb × workflow.py × serving.py × secret.txt × load_data.py × trainer.py × +
1 import mlrun
2 from sklearn.ensemble import RandomForestClassifier
3 from sklearn.model_selection import train_test_split
4 from mlrun.frameworks.sklearn import apply_mlrun
5
6 def train(
7     dataset: mlrun.DataItem,
8     label_column: str = 'target',
9     n_estimators: int = 100,
10    max_depth: int = None,
11    model_name: str = "breast_cancer_classifier"
12):
13    df = dataset.as_df()
14    X = df.drop(label_column, axis=1)
15    y = df[label_column]
16    X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.1, random_state=42)
17    model = RandomForestClassifier(n_estimators=n_estimators, max_depth=max_depth, random_state=42)
18    apply_mlrun(model=model, model_name=model_name, x_test=X_test, y_test=y_test)
19    model.fit(X_train, y_train)
20
```

Serving.py

The Classifier Model inherits from **mlrun.serving.V2ModelServer**. Thus, we can load the model from storage and be able to predict according to requests.



```
1 from cloudpickle import load
2 import numpy as np
3 from typing import List
4 import mlrun
5
6 class ClassifierModel(mlrun.serving.V2ModelServer):
7     def load(self):
8
9         model_file, extra_data = self.get_model('.pkl')
10        self.model = load(open(model_file, 'rb'))
11
12    def predict(self, body: dict) -> List:
13
14        feats = np.asarray(body['inputs'])
15        results: np.ndarray = self.model.predict(feats)
16        return results.tolist()
```

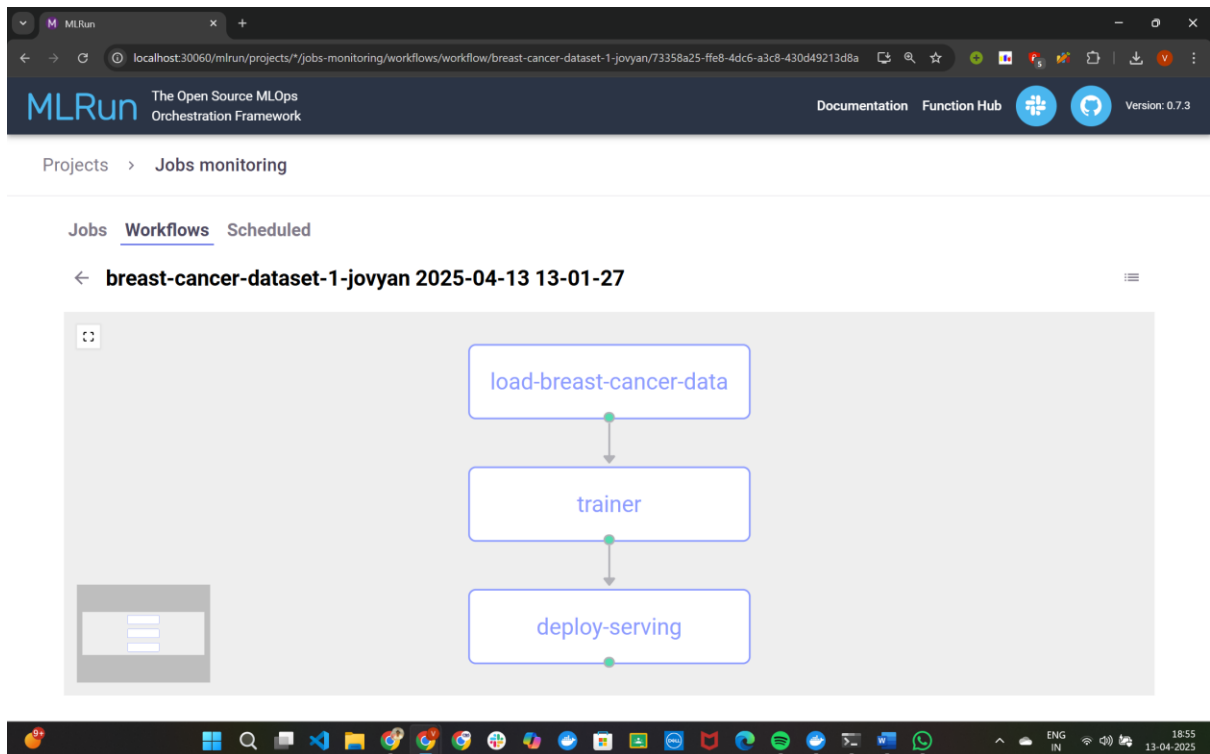
Workflow.py

Data Ingestion:

The **@dsl.pipeline** decorator. For data ingestion, the **load-data** function from **load_data.py** is ran and the output is stored in **dataset**. The **trainer** function in **trainer.py** is used to experiment with different hyperparameters like number of estimators, and max depth. The **max_accuracy** selector chooses the model with the best accuracy, effectively performing a grid search over the hyperparameters. The output is stored in **model**. Next, the model is deployed using the ClassifierModel in **serving.py**.

```
main.ipynb x workflow.py x serving.py x project.yaml x secret.txt x load_data.py x trainer.py x
1 import mlrun
2 from kfp import dsl
3
4 @dsl.pipeline(name="breast-cancer-pipeline")
5 def pipeline(model_name="breast_cancer_classifier"):
6
7     ingest = mlrun.run_function(
8         "load-data",
9         name="load-breast-cancer-data",
10         params={"format": "pq", "model_name": model_name},
11         outputs=["dataset"],
12     )
13
14     train = mlrun.run_function(
15         "trainer",
16         inputs={"dataset": ingest.outputs["dataset"]},
17         hyperparams={
18             "n_estimators": [10, 100, 200],
19             "max_depth": [2, 5, 10]
20         },
21         selector="max.accuracy",
22         outputs=["model"],
23     )
24
25     deploy = mlrun.deploy_function(
26         "serving",
27         models=[{"key": model_name, "model_path": train.outputs["model"], "class_name": "ClassifierModel"}],
28         mock=False
29     )
30
31
```

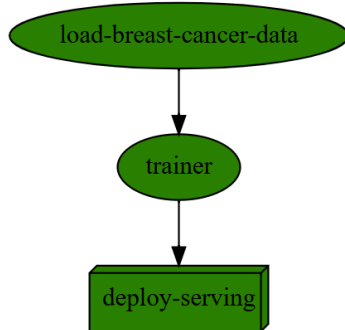
Workflow Graph



[click here to view progress](#)

Pipeline running (id=73358a25-ffe8-4dc6-a3c8-430d49213d8a), [click here](#) to view the details in MLRun UI

```
> 2025-04-13 13:01:28,393 [info] Started run workflow breast-cancer-dataset-1-jovyan with run id = '73358a25-ffe8-4dc6-a3c8-430d49213d8a' by kfp engine
> 2025-04-13 13:01:28,394 [info] Waiting for pipeline run completion: {"project": "breast-cancer-dataset-1-jovyan", "run_id": "73358a25-ffe8-4dc6-a3c8-430d49213d8a"}
```



click the hyper links below to see detailed results

uid	start	state	kind	name	parameters	results
...	Apr 13 13:02:12	completed	run	trainer		best_iteration=1 accuracy=0.9649122807017544 f1_score=0.975 precision_score=0.975 recall_score=0.975
...	Apr 13 13:01:48	completed	run	load-breast-cancer-data	format=pq model_name=breast-cancer-classifier	label_column=target

Data_prep artifact

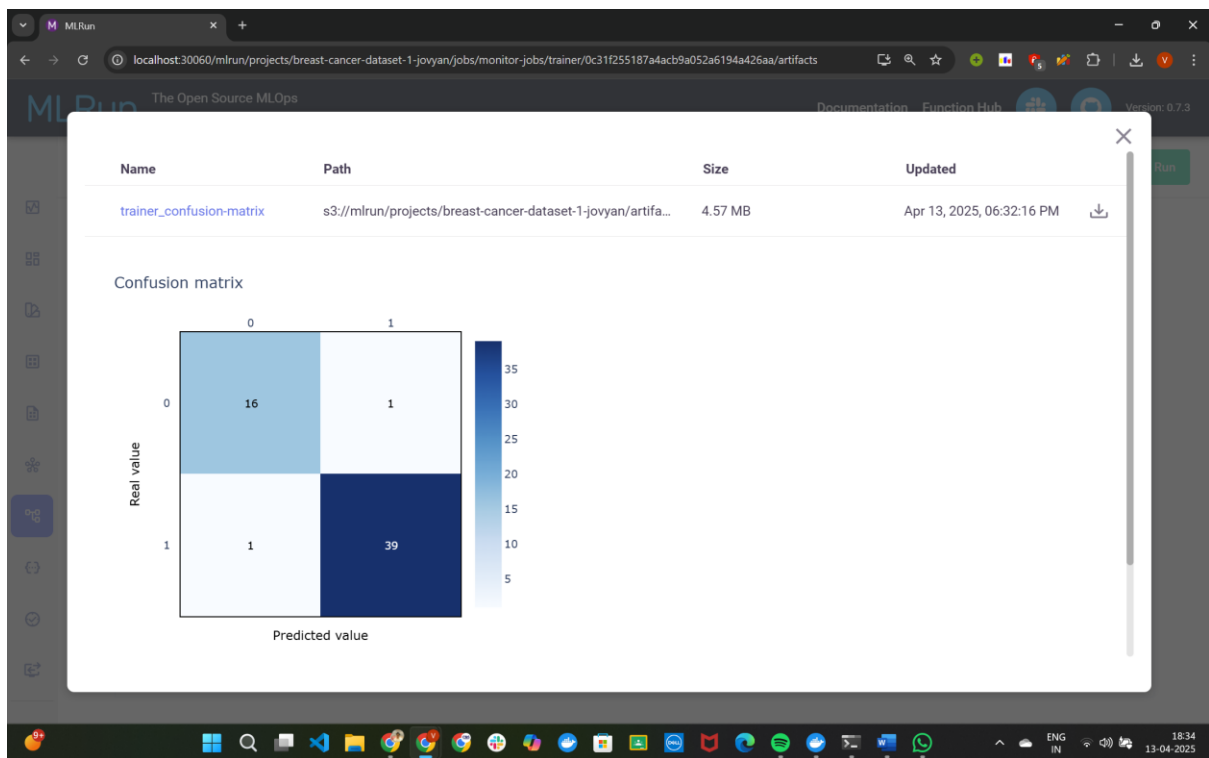
The screenshot displays the MLRun UI for monitoring a job. The main view is titled 'load-breast-cancer-data' and shows the 'Artifacts' tab. A table lists the artifact details:

Name	Path	Size	Updated
load-breast-cancer-data_breas...	s3://mlrun/projects/breast-cancer-dataset-1-jovyan/artifacts/7...	148 kB	Apr 13, 2025, 06:31:...

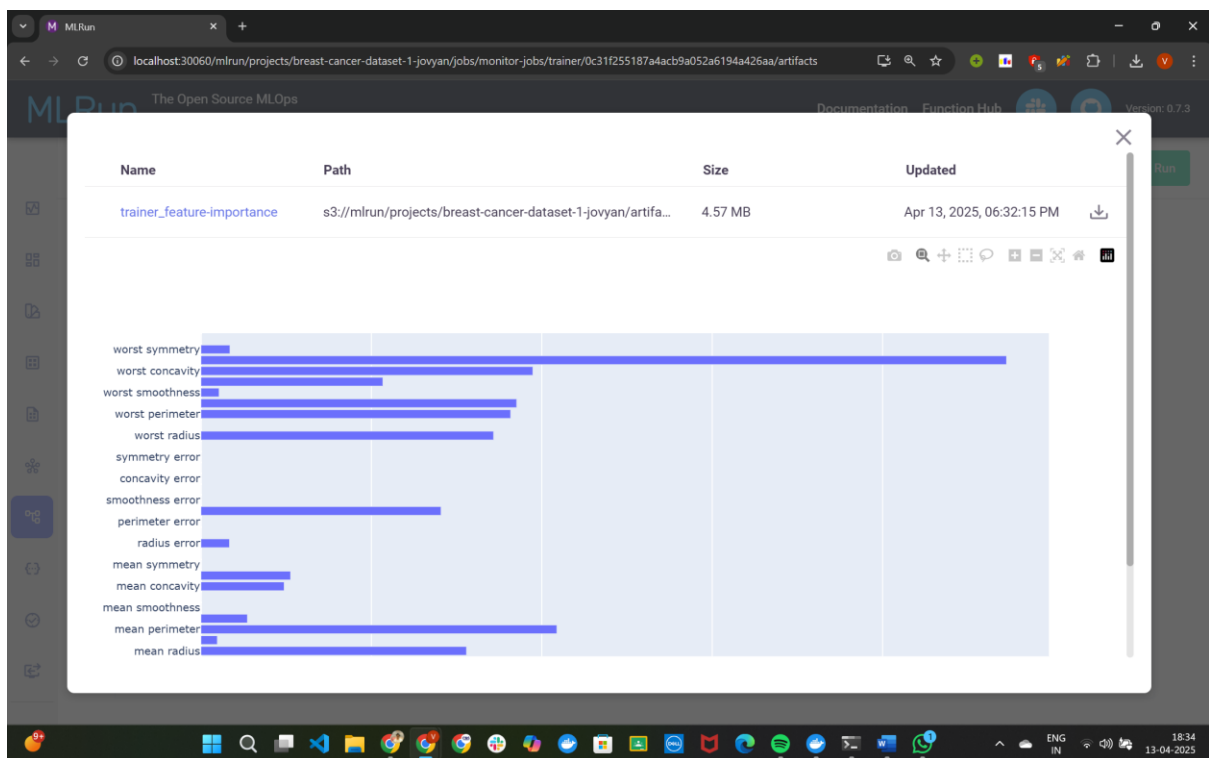
Below the table, a data preview is shown with columns: index, mean radius, mean texture, mean perimeter, mean area, mean smooth..., mean compac..., mean concavity, and mean conc. The preview shows two rows of data.

index	mean radius	mean texture	mean perimeter	mean area	mean smooth...	mean compac...	mean concavity	mean concav...	mean symmet...
9	12.46	24.04	83.97	475.9	0.1186	0.2396	0.2273	0.08543	0.203
10	16.02	23.24	102.7	797.8	0.08206	0.06669	0.03299	0.03323	0.1528
11	15.78	17.89	103.6	781	0.0971	0.1292	0.09954	0.06606	0.1842
12	19.17	24.8	132.4	1123	0.0974	0.2458	0.2065	0.1118	0.2397
13	15.85	23.95	103.7	782.7	0.08401	0.1002	0.09938	0.05364	0.1847
14	13.73	22.61	93.6	578.3	0.1131	0.2293	0.2128	0.08025	0.2069
15	14.54	27.54	96.73	658.8	0.1139	0.1595	0.1639	0.07364	0.2303
16	14.68	20.13	94.74	684.5	0.09867	0.072	0.07395	0.05259	0.1586
17	16.13	20.68	108.1	798.8	0.117	0.2022	0.1722	0.1028	0.2164
18	19.81	22.15	130	1260	0.09831	0.1027	0.1479	0.09498	0.1582

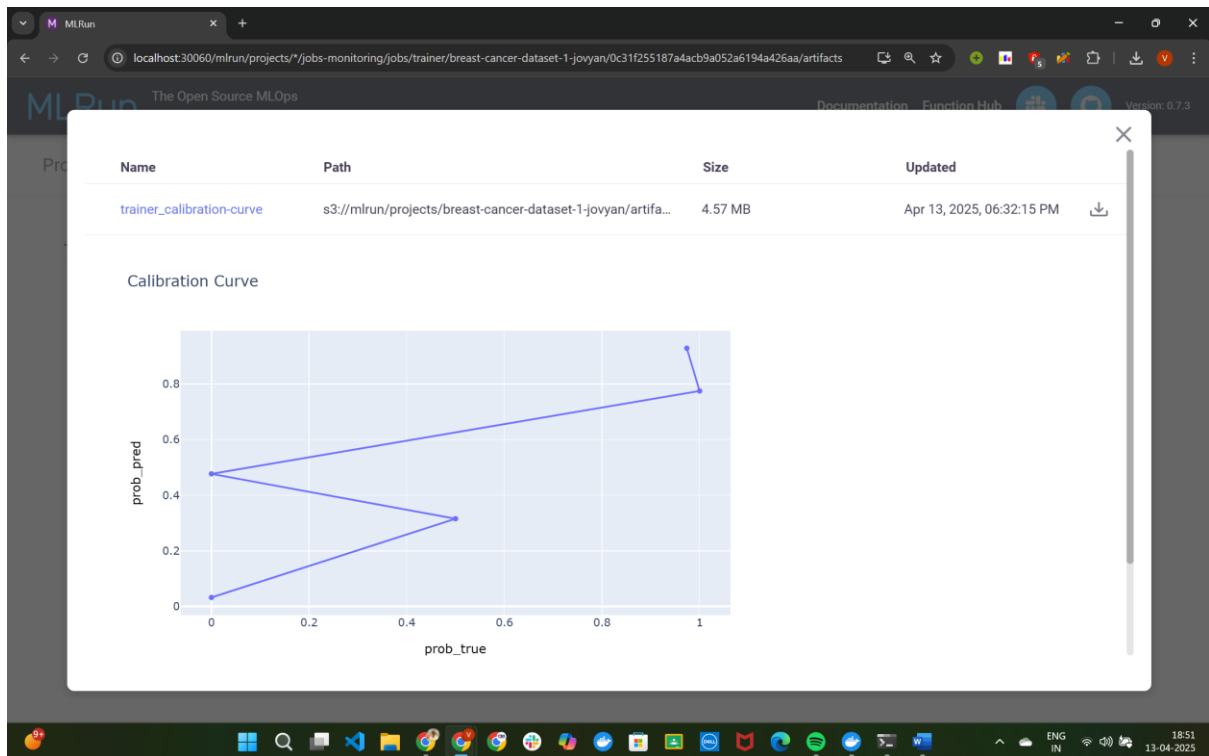
Confusion Matrix



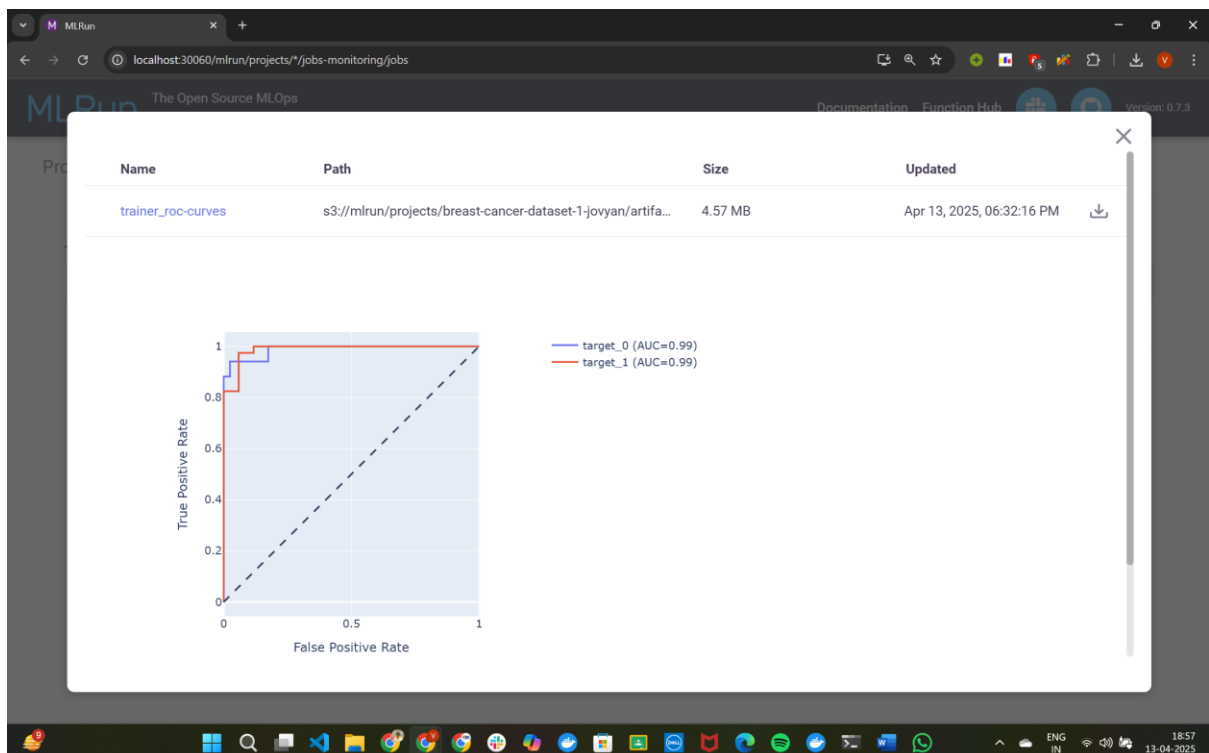
Feature Selection Artifact



Calibration Curve



Trainer – ROC Curves



Project.yaml

```
main.ipynb x workflow.py x serving.py x secret.txt x project.yaml x load_data.py x trainer.py x +
1 metadata:
2   name: breast-cancer-dataset-1-jovyan
3 spec:
4   functions:
5     - url: load_data.py
6       name: load-data
7       kind: job
8       image: mlrun/mlrun
9     - url: trainer.py
10      name: trainer
11      kind: job
12      image: mlrun/mlrun
13      handler: train
14     - url: serving.py
15       name: serving
16       kind: serving
17       image: mlrun/mlrun
18   origin_url: git://Venkat-2341:ghp_myqkaWfNrBt68hkCdIB4fJHwpmi72d3MSgrF@github.com/Venkat-2341/lab9.git#refs/heads/master
19   conda: ''
20   source: git://Venkat-2341:ghp_myqkaWfNrBt68hkCdIB4fJHwpmi72d3MSgrF@github.com/Venkat-2341/lab9.git#refs/heads/master
21   desired_state: online
22 kind: project
23
```

MLRun The Open Source MLOps Orchestration Framework Documentation Function Hub Version: 0.7.3

Projects > breast-cancer-dataset-1-jovyan > Jobs and workflows Batch Run

Monitor Jobs Monitor Workflows Schedule

← **trainer** Iteration: 1 (Best iteration) Resource monitoring

Apr 13, 2025, 06:32:42 PM

Overview Inputs Artifacts Results Logs Pods

Name ↑	Path	Size	Updated	
breast_cancer_classifier	s3://mlrun/projects/breast-cancer-dataset-1-jovyan/artifacts/7...	10.6 kB	Apr 13, 2025, 06:32:...	
trainer_calibration-curve	s3://mlrun/projects/breast-cancer-dataset-1-jovyan/artifacts/7...	4.57 MB	Apr 13, 2025, 06:32:...	
trainer_confusion-matrix	s3://mlrun/projects/breast-cancer-dataset-1-jovyan/artifacts/7...	4.57 MB	Apr 13, 2025, 06:32:...	
trainer_feature-importance	s3://mlrun/projects/breast-cancer-dataset-1-jovyan/artifacts/7...	4.57 MB	Apr 13, 2025, 06:32:...	
trainer_roc-curves	s3://mlrun/projects/breast-cancer-dataset-1-jovyan/artifacts/7...	4.57 MB	Apr 13, 2025, 06:32:...	