

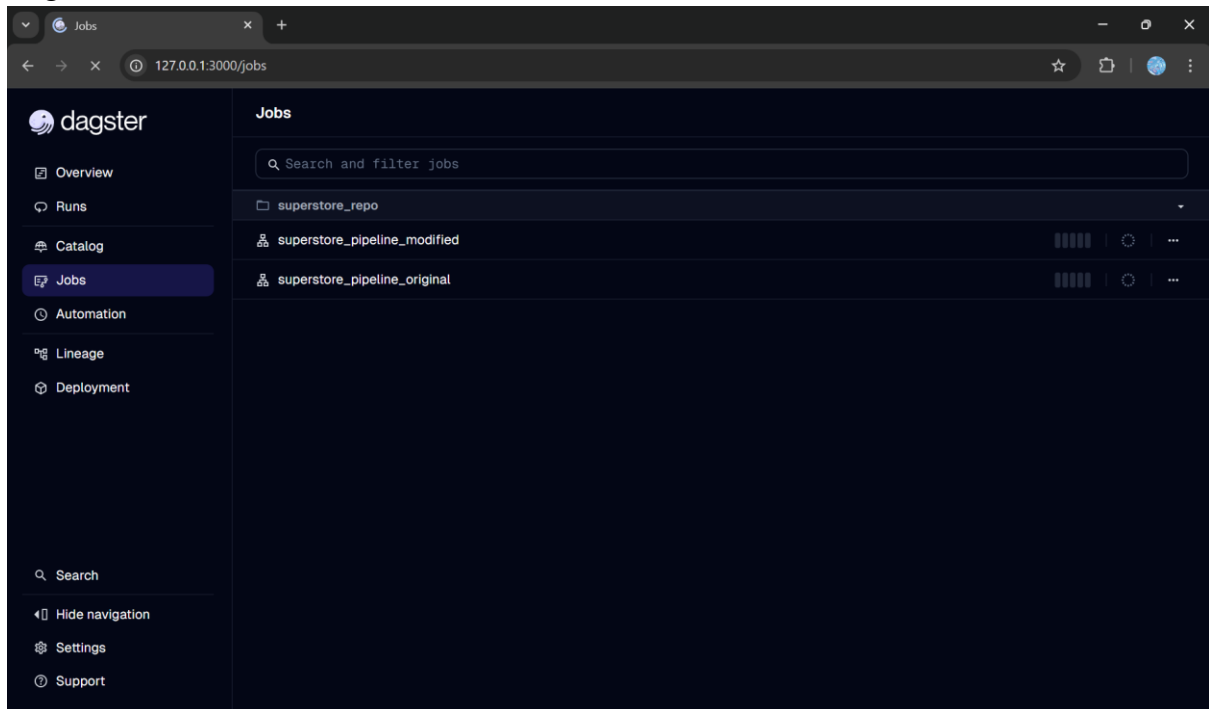
Report: Time Savings and Reproducibility Using Dagster ML Pipeline

Conventional machine learning workflows run in Jupyter notebooks can struggle with reproducibility because of the non-linear execution of cells, modifications of code, and changes made to the versions of data. Often, these problems lead to inconsistent results and hinder effective tracking of experiments. To mitigate these challenges, this project investigates the use of Dagster as a data orchestration framework for building a structured and reproducible machine learning pipeline.

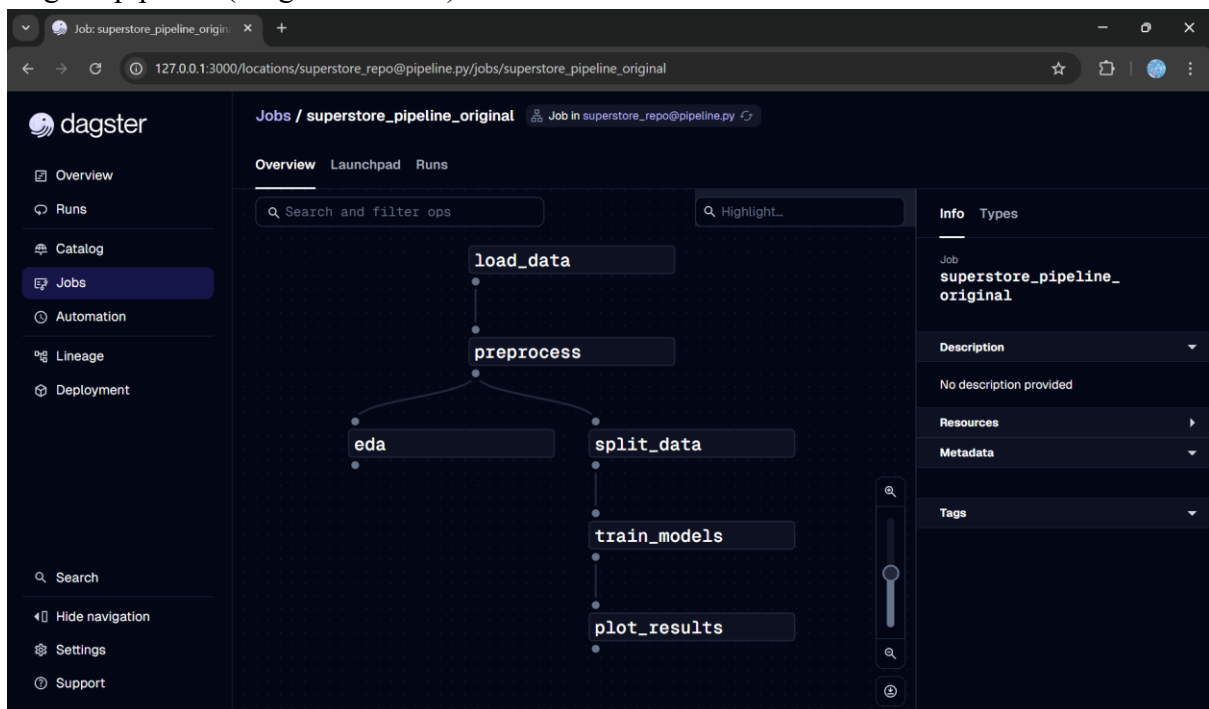
The pipeline was built using the Superstore dataset and goes through all the primary steps of a standard machine learning workflow: data ingestion and preprocessing, model training and evaluation, and visualization of results. Three different types of classification models (Decision Tree, Random Forest, and Logistic Regression) were built and evaluated. Dagster guarantees that there are explicit dependencies established between different stages of the pipeline, which will ensure consistent execution orders and reproducible data flows.

Reproducibility of both runs was demonstrated by running the same pipeline against the original and a modified dataset without changing code in either version of the pipeline. While both runs of the pipelines were tracked by Dagster as different runs, the timing difference between the two pipelines was measurable. While the workflow was executed in a traditional Jupyter Notebook, the original dataset took **2.80 seconds** to execute, while the modified dataset required only **2.33 seconds** to complete; in contrast, the Dagster-based pipeline took **1 minute 14 seconds** for the original dataset and **1 minute 12 seconds** for the modified dataset. While the notebook-based approach was substantially faster for this dataset size than the Dagster-based pipeline, there are clear advantages to using Dagster with respect to execution control, reproducibility, and workflow organization, all of which make it a better fit for machine learning pipelines scaled and deployed into production environments where reliability and traceability are critically important.

Dagster:



Dagster pipeline (Original Dataset):



The top screenshot shows the Dagster 'Launchpad' for the job 'superstore_pipeline_original'. The interface includes a sidebar with navigation links (Overview, Runs, Catalog, Jobs, Automation, Lineage, Deployment) and a main area with a 'New Run' button and a code editor. The code editor contains a pipeline configuration snippet:

```

/* Configure how steps are executed
within a run. */
execution?: {
  config?: ...
}
/* Configure how loggers emit
messages within a run. */
loggers?: {
  console?: ...
}
/* Configure runtime parameters for
ops or assets. */
ops?: {
  eda?: ...
  load_data?: ...
  plot_results?: ...
  preprocess?: ...
  split_data?: ...
}

```

The bottom screenshot shows the 'Runs' view for a specific run (46a3c847). It displays a progress bar for the 'eda' step and a table of events:

TIMESTAMP	OP	EVENT TYPE	INFO
12:37:12.572 pm	plot_results	STEP_OUTPUT	Yielded output "result" of type "Any". (Type check passed).
12:37:12.642 pm	plot_results	HANDLED_OUTPUT	Handled output "result" using IO manager "io_manager"
12:37:12.665 pm	plot_results	STEP_SUCCESS	Finished execution of step "plot_results" in 2.0s.
12:37:15.771 pm	-	ENGINE_EVENT	Multiprocess executor: parent process exiting after 1m13s (pid: 32784)
12:37:15.804 pm	-	RUN_SUCCESS	Finished execution of run for "superstore_pipeline_original".
12:37:15.966 pm	-	ENGINE_EVENT	Process for run exited (pid: 32784).

Dagster pipeline (Modified Dataset):

dagster

Overview

Runs

Catalog

Jobs

Automation

Lineage

Deployment

Search

Hide navigation

Settings

Support

Jobs / superstore_pipeline_modified

Job in superstore_repo@pipeline.py

OverviewLaunchpadRuns

Search and filter ops

Highlight...

load_data_modified

preprocess

eda

split_data

train_models

plot_results

Info

Types

Job

superstore_pipeline_modified

Description

No description provided

Resources

Metadata

Tags

Jobs / superstore_pipeline_modified

Job in superstore_repo@pipeline.py

OverviewLaunchpadRuns

New Run

+ Add...

✖ *

Edit tags

▶

1 {}

2

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

execution?: {

config?: ...

}

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

loggers?: {

console?: ...

}

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

ops?: {

eda?: ...

load_data_modified?: ...

plot_results?: ...

}

Use Ctrl+Space to show auto-completions inline.

ERRORS

No errors

CONFIG ACTIONS:

Scaffold missing config

No missing config

RUNTIME

executionloggersio_manager

OPS

eda

load_data_modified

plot_results

preprocess

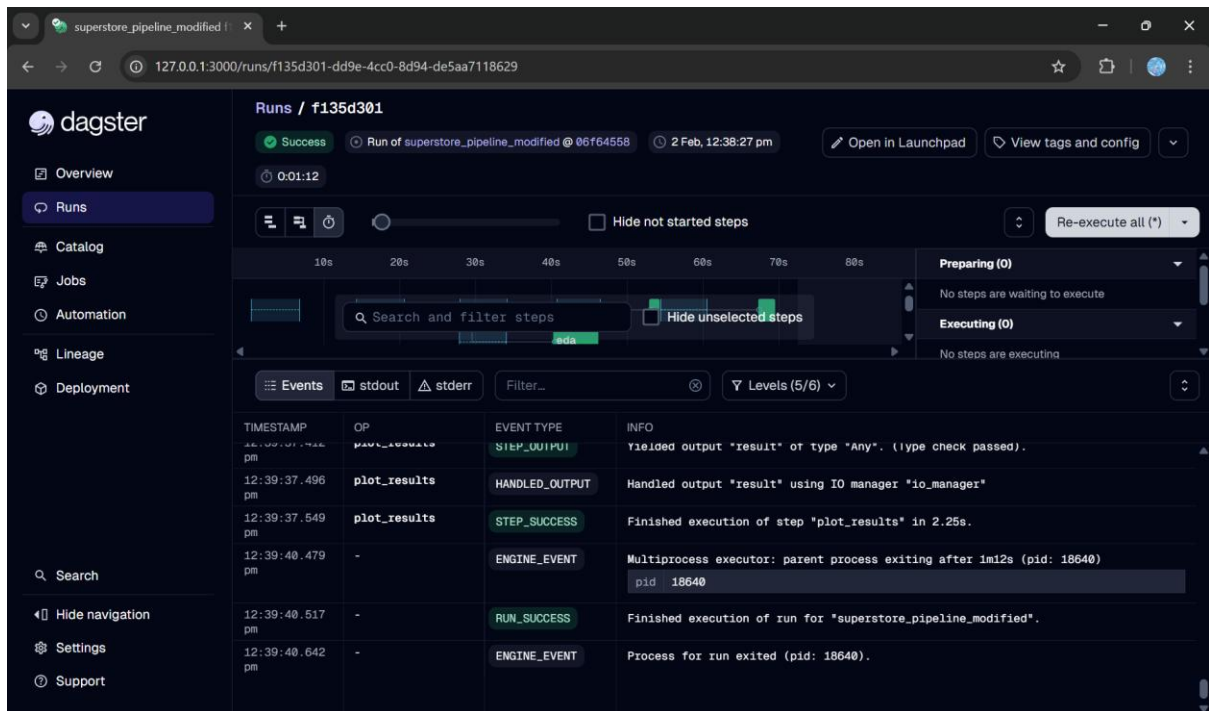
split_data

Launch Run

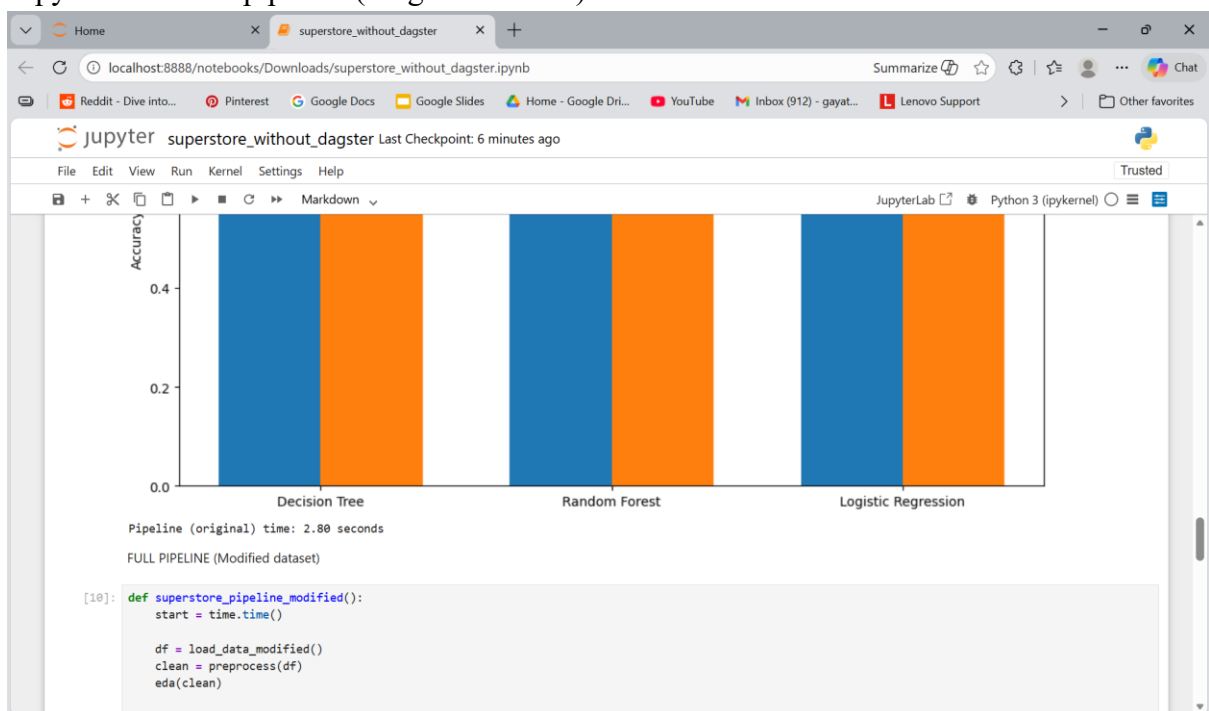
RESOURCES

Errors only

4



Jupyter Notebook pipeline (Original Dataset):



Jupyter Notebook pipeline (Modified Dataset):

