

NAME : NAVIYA DHARSHINI A S
ROLL NO : 23AD083
DAY -08 - 11/07/2025

DVC (Data Version Control) - Full Notes from Basics to Advanced



1. What is DVC?

- DVC = Git for Data + ML Pipelines
 - Manages **data, models, and experiments** in ML projects.
 - Works with Git but tracks **large files, directories, and ML stages**.
 - Makes projects **reproducible, shareable, and collaborative**.
-



2. Why DVC?



Git Problem

Can't track large datasets

DVC tracks large files via `.dvc` metadata



DVC Solution

Git lacks pipeline tracking

DVC handles pipelines + stages like `make`

No reproducibility for ML

DVC ensures same input = same output

Can't share data easily

DVC supports remote storage (S3, GDrive etc)



4. DVC Project Structure

```
my-ml-project/
├── data/                      # Raw datasets (not tracked by Git)
├── data.dvc                   # DVC metadata for dataset
├── model.pkl                  # Trained model (tracked by DVC)
├── model.pkl.dvc
├── dvc.yaml                  # ML pipeline stages
├── dvc.lock                   # Exact command runs (hashes etc.)
├── .dvc/
└── .git/                      # Git repo
```

Screenshot of the AWS IAM 'Create policy' wizard, Step 1: Specify permissions.

The Policy editor shows the following JSON:

```

1  {
2    "Version": "2012-10-17",
3    "Statement": [
4      {
5        "Effect": "Allow",
6        "Action": [
7          "s3:ListBucket",
8          "s3:GetObject",
9          "s3:PutObject",
10         "s3:DeleteObject",
11         "s3:HeadBucket"
12       ],
13       "Resource": [
14         "arn:aws:s3:::dvcbucketml",
15         "arn:aws:s3:::dvcbucketml*"
16     ]
}

```

The 'Add actions' sidebar shows services like AI Operations, AMP, API Gateway, and API Gateway V2.

Screenshot of the AWS IAM 'Create policy' wizard, Step 2: Review and create.

The 'Permissions defined in this policy' section shows the following table:

Allow (1 of 445 services)			
Service	Access level	Resource	Request condition
S3	Limited: List, Read, Write	Multiple	None

The 'Add tags - optional' section shows 'No tags associated with the resource'.

Buttons at the bottom: Cancel, Previous, Create policy.

Specify user details

User details

User name: dvuser

The user name can have up to 64 characters. Valid characters: A-Z, a-z, 0-9, and + = . @ _ - (hyphen)

Provide user access to the AWS Management Console - *optional*
If you're providing console access to a person, it's a [best practice](#) to manage their access in IAM Identity Center.

If you are creating programmatic access through access keys or service-specific credentials for AWS CodeCommit or Amazon Keyspaces, you can generate them after you create this IAM user. [Learn more](#)

Cancel **Next**

Add user to an existing group or create a new one. Using groups is a best-practice way to manage user's permissions by job functions. [Learn more](#)

Add user to group
Add user to an existing group, or create a new group. We recommend using groups to manage user permissions by job function.

Copy permissions
Copy all group memberships, attached managed policies, and inline policies from an existing user.

Attach policies directly
Attach a managed policy directly to a user. As a best practice, we recommend attaching policies to a group instead. Then, add the user to the appropriate group.

Create policy

Permissions policies (1/1379)
Choose one or more policies to attach to your new user.

Filter by Type: Q: dv, All types, 1 match

Policy name	Type	Attached entities
<input checked="" type="checkbox"/> dvpolicy	Customer managed	0

Set permissions boundary - *optional*

Cancel **Previous** **Next**

Review your choices. After you create the user, you can view and download the autogenerated password, if enabled.

User details

User name: dvuser	Console password type: None	Require password reset: No
-------------------	-----------------------------	----------------------------

Permissions summary

Name	Type	Used as
<input checked="" type="checkbox"/> dvpolicy	Customer managed	Permissions policy

Tags - *optional*
Tags are key-value pairs you can add to AWS resources to help identify, organize, or search for resources. Choose any tags you want to associate with this user.
No tags associated with the resource.

Add new tag
You can add up to 50 more tags.

Cancel **Previous** **Create user**

Screenshot of the AWS IAM User Details page for 'dvcuser'.

Identity and Access Management (IAM)

ARN: arn:aws:iam::929690601290:user/dvcuser

Created: July 11, 2025, 10:35 (UTC+05:30)

Console access: Disabled

Last console sign-in: -

Access key 1: Create access key

Permissions: Groups, Tags, Security credentials, Last Accessed

Permissions policies (1): dvcpolicy (Customer managed, Directly attached)

Permissions boundary (not set):

Generate policy based on CloudTrail events: You can generate a new policy based on the access activity for this user, then customize, create, and attach it to this role. AWS uses your CloudTrail events to identify the services and actions.

Screenshot of the 'Create access key' wizard step 2: 'Set description tag - optional'.

Local code: You plan to use this access key to enable application code in a local development environment to access your AWS account.

Application running on an AWS compute service: You plan to use this access key to enable application code running on an AWS compute service like Amazon EC2, Amazon ECS, or AWS Lambda to access your AWS account.

Third-party service: You plan to use this access key to enable access for a third-party application or service that monitors or manages your AWS resources. (Selected)

Application running outside AWS: You plan to use this access key to authenticate workloads running in your data center or other infrastructure outside of AWS that needs to access your AWS resources.

Other: Your use case is not listed here.

Alternative recommended: As a best practice, use temporary security credentials (IAM roles) instead of creating long-term credentials like access keys, and don't create AWS account root user access keys. [Learn more](#)

Confirmation: I understand the above recommendation and want to proceed to create an access key.

Next

Screenshot of the 'Create access key' wizard step 3: 'Set description tag - optional'.

Step 1: Access key best practices & alternatives

Step 2 - optional: Set description tag (Selected)

Step 3: Retrieve access keys

Set description tag - optional: Info

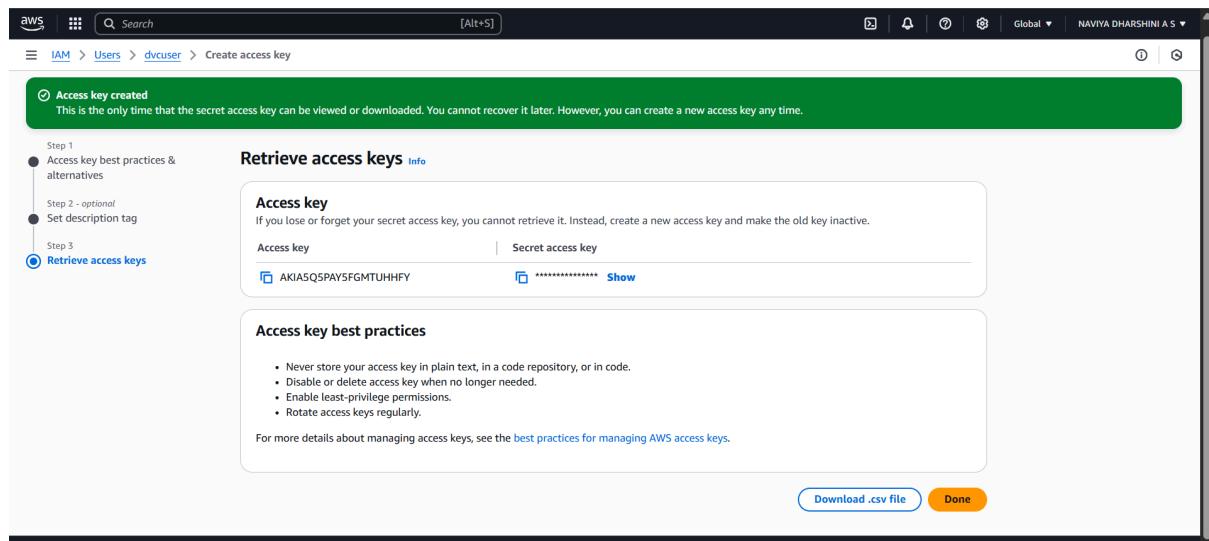
The description for this access key will be attached to this user as a tag and shown alongside the access key.

Description tag value: dvctagvalue

Describe the purpose of this access key and where it will be used. A good description will help you rotate this access key confidently later.

Maximum 256 characters. Allowed characters are letters, numbers, spaces representable in UTF-8, and: _ . / = + - @

Cancel **Previous** **Create access key**



🌐 . DVC Remotes (Cloud Backends)

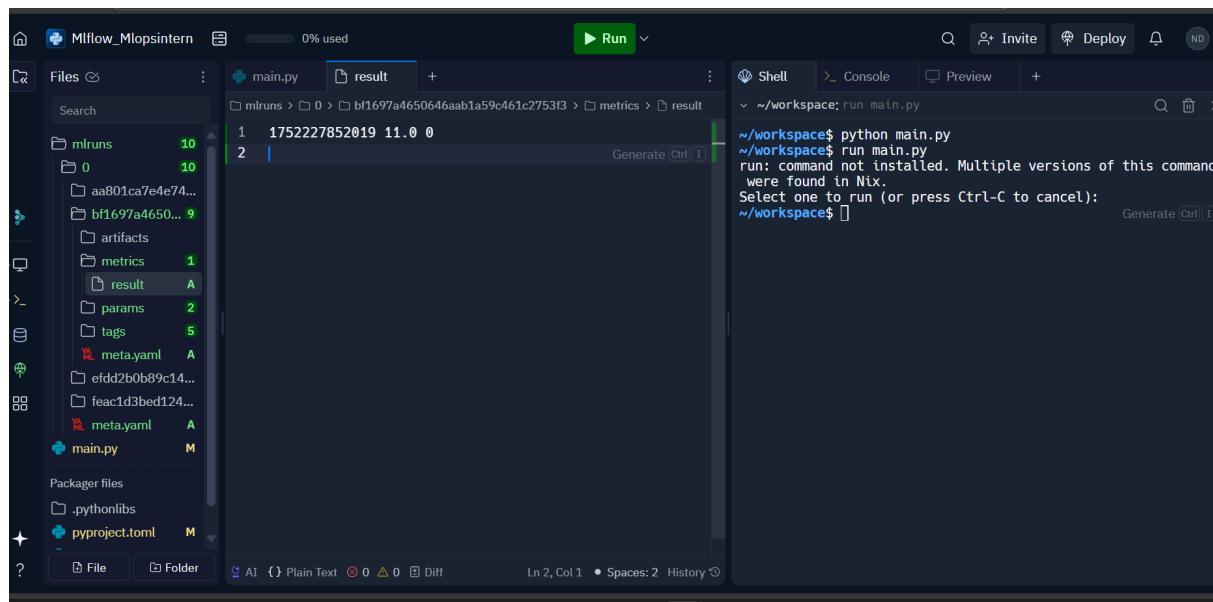
Supports:

- ✓ AWS S3
- ✓ Google Drive
- ✓ Azure Blob
- ✓ GCP Bucket
- ✓ SSH, HDFS, WebDAV

```
bash
CopyEdit
dvc remote add -d storage s3://mybucket
dvc remote modify storage access_key_id ...
dvc push
```

MLflow is a platform to manage the **machine learning lifecycle**. It helps **track, compare, and reproduce** ML experiments efficiently. Here's how the **end-to-end flow** works, from start to finish:

USING REPLIT :



Files

main.py result

mlruns 10

0 10

aa801ca7e4e74...
bf1697a4650... 9

artifacts
metrics 1
result A
params 2
tags 5
meta.yaml A
meta.yaml A

main.py M

.pythonlibs
pyproject.toml M

Shell

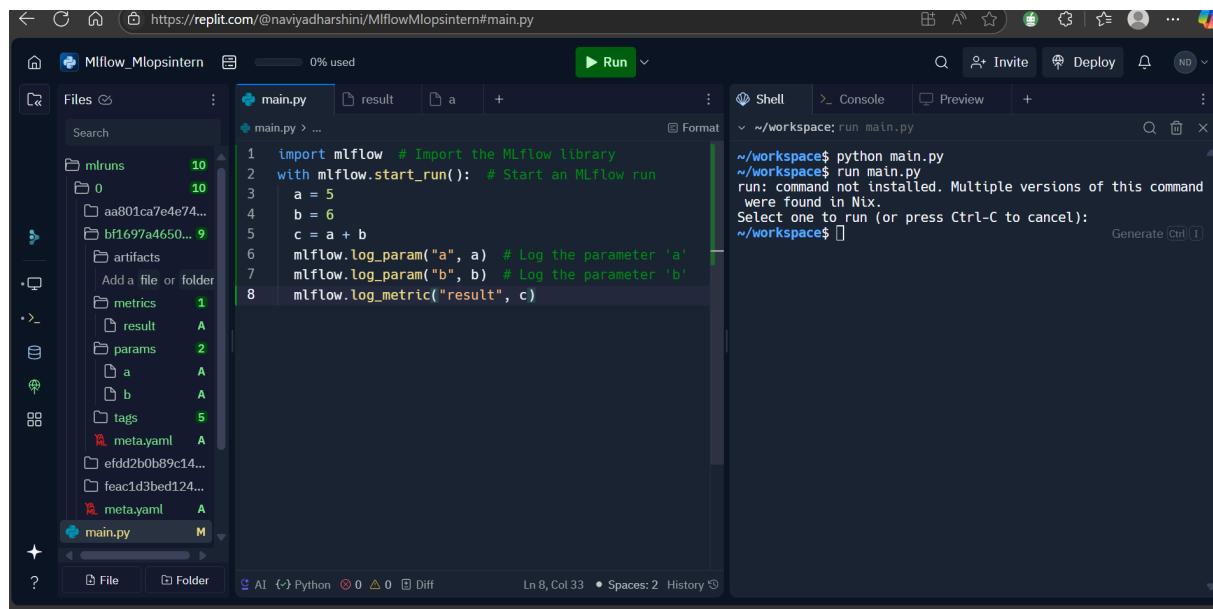
~/workspace\$ python main.py

~/workspace\$ run main.py

run: command not installed. Multiple versions of this command were found in Nix.

Select one to run (or press Ctrl-C to cancel):

~/workspace\$



https://replit.com/@naviyadharshini/MLflowMlopsintern#main.py

MLflow_Mlopsintern 0% used

Files

- mlruns (10)
- 0 (10)
- aa801ca7e4e74...
- bf1697a4650... (9)
- artifacts
- metrics (1)
- result (A)
- params (2)
- a (A)
- b (A)
- tags (5)
- meta.yaml (A)
- efdd2b0b89c14...
- feac1d3bed124...
- meta.yaml (A)

main.py

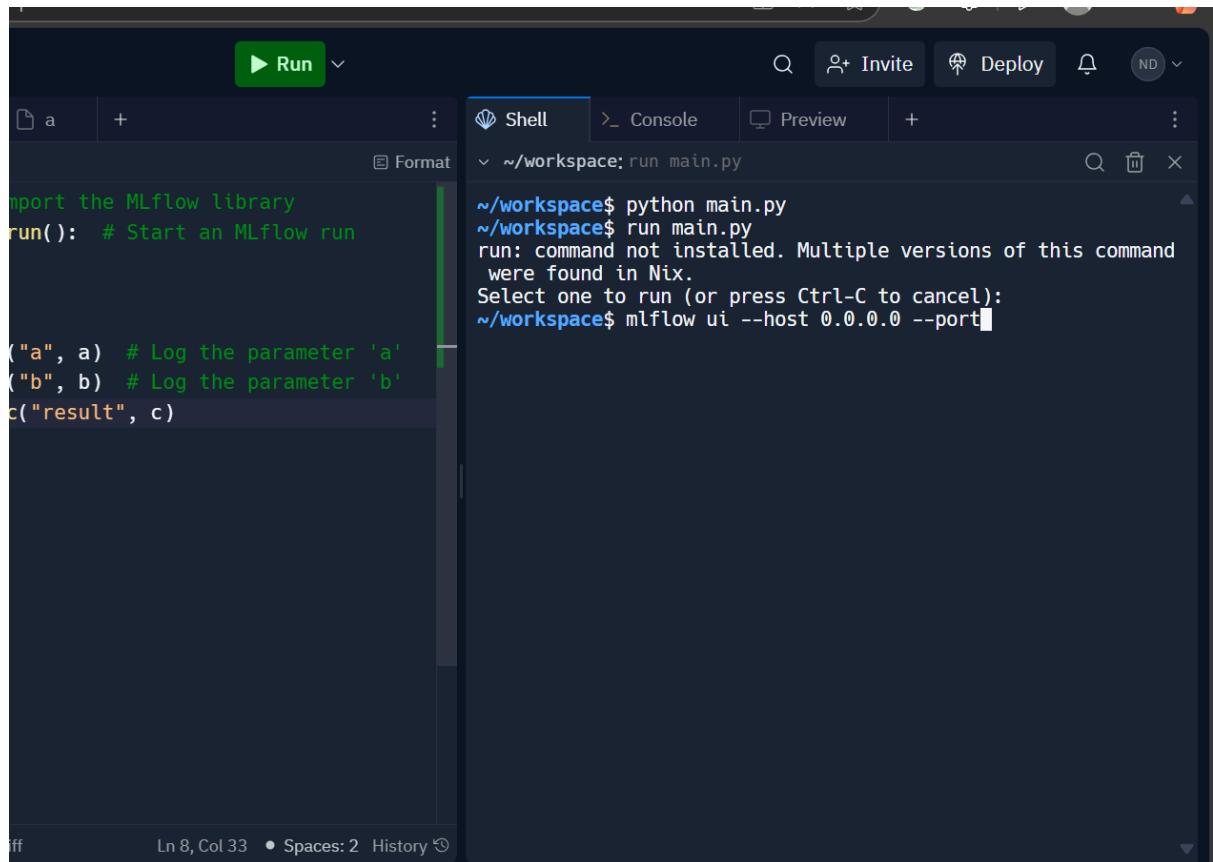
```
1 import mlflow # Import the MLflow library
2 with mlflow.start_run(): # Start an MLflow run
3     a = 5
4     b = 6
5     c = a + b
6     mlflow.log_param("a", a) # Log the parameter 'a'
7     mlflow.log_param("b", b) # Log the parameter 'b'
8     mlflow.log_metric("result", c)
```

Run

Shell

```
~/workspace$ python main.py
~/workspace$ run main.py
run: command not installed. Multiple versions of this command
were found in Nix.
Select one to run (or press Ctrl-C to cancel):
~/workspace$
```

Invite Deploy ND



Run

Files

- a
- +

main.py

```
import the MLflow library
run(): # Start an MLflow run

("a", a) # Log the parameter 'a'
("b", b) # Log the parameter 'b'
c("result", c)
```

Run

Shell

```
~/workspace$ python main.py
~/workspace$ run main.py
run: command not installed. Multiple versions of this command
were found in Nix.
Select one to run (or press Ctrl-C to cancel):
~/workspace$ mlflow ui --host 0.0.0.0 --port
```

Invite Deploy ND

MLflow UI (version 1.30.0) - Mlflow_Mlopsintern

0% used

Run: **Run** ▾

Shell ▾ Console Preview +

Experiments Default ▾ Provide Feedback Add Description Share

Runs Models Experimental Evaluation Traces

Search experiments: metrics.rmse < 1 and params.model = "tree"

Time created ▾ State: Active ▾ Datasets ▾ Sort: Created ▾

Columns ▾ Group by ▾

Run Name	Created	Dataset
illustrious-fox-469	1 minute ago	-
learned-fawn-378	3 minutes ago	-
adventurous-worm-912	3 minutes ago	-
youthful-kit-992	3 minutes ago	-

5 matching runs

File Folder

MLflow UI (version 3.1.1) - Mlflow_Mlopsintern

0% used

Run: **Run** ▾

Shell ▾ Console Preview +

mlflow 3.1.1 Experiments Models Prompts GitHub Docs

Experiments Default ▾ Provide Feedback Add Description Share

Runs Models Experimental Evaluation Traces

Rename Delete Compare Add tags ▾

Run Name	Created	Dataset
illustrious-fox-469	2 minutes ago	-
learned-fawn-378	4 minutes ago	-
adventurous-worm-912	4 minutes ago	-
youthful-kit-992	4 minutes ago	-

5 matching runs

File Folder

MLflow UI (version 3.1.1) - Mlflow_Mlopsintern

0% used

Run: **Run** ▾

Shell ▾ Console Preview +

Comparing 2 Runs from 1 Experiment

Visualizations

Parallel Coordinates Plot Scatter Plot Box Plot Contour Plot

Parameters: Please select parameters

Metrics: result ▾

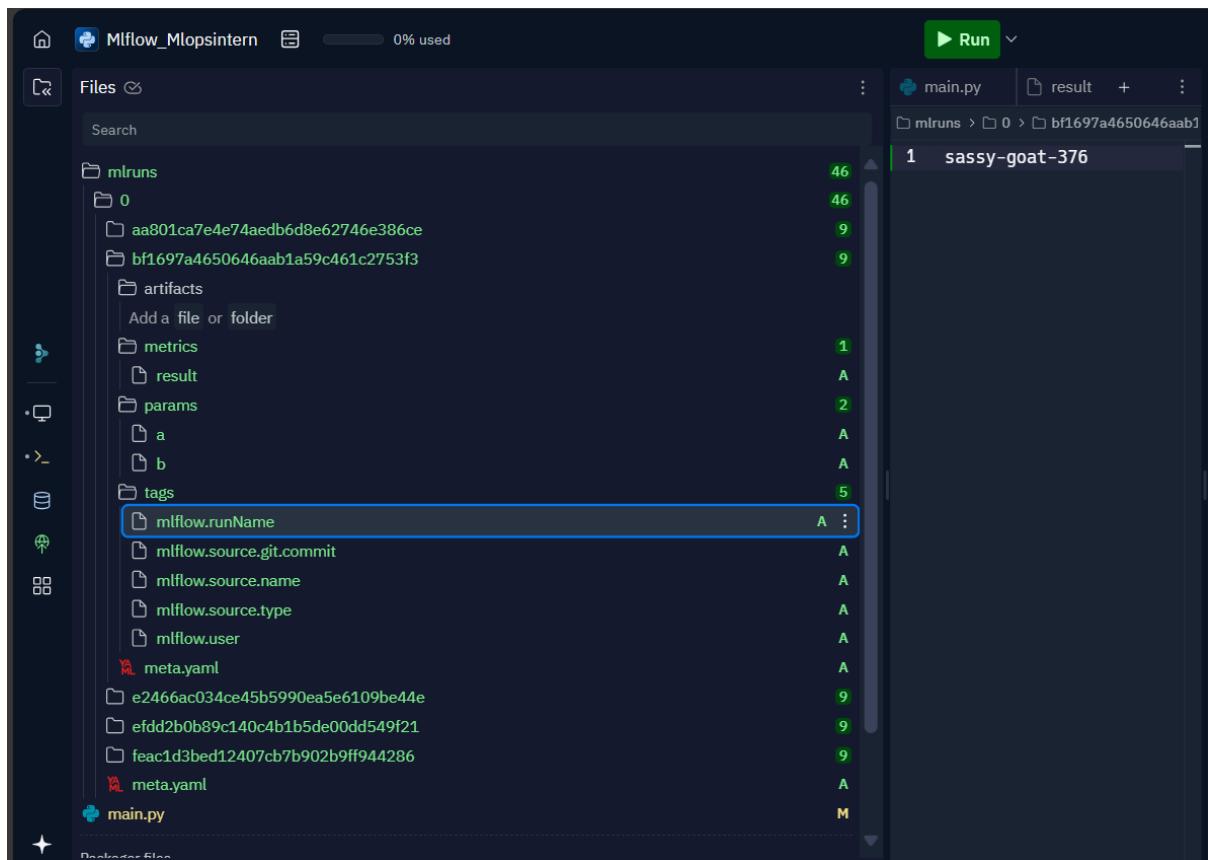
Clear All

result

12.10000
12.00000
11.50000
11.00000
10.50000

11.4
11.2
11
10.8
10.6

File Folder



1 Experiment

- ◆ Think of an **Experiment** as a **folder** or **container** for all ML model tests.
- ◆ It groups multiple **runs** under one goal — like testing different algorithms or hyperparameters to solve a problem (e.g., predicting student results or sales forecast).

 Example:

“Student Grade Prediction” can be the experiment name.

2 Run

- ◆ A **Run** is a **single trial or training session** within an experiment.
- ◆ Every time you train a model (with different inputs or settings), it's recorded as one run.

 Purpose: Track the details of that attempt — like which model was used, how it performed, etc.

Example: Trying a Random Forest with 100 trees = 1 run
Changing it to 200 trees = another run

3 Parameters

- These are the **inputs** or **settings** used in each run. Think of them like the ingredients in a recipe.
- They include things like learning rate, number of estimators, batch size, etc.

 Benefit: You can compare which parameters gave the best results.

Example: `n_estimators=100, max_depth=10`

4 Metrics

- These are the **outcomes** or **performance scores** measured during training/testing.
- Includes metrics like accuracy, F1 score, precision, loss, etc.

 Goal: Helps decide which run/model performs best.

Example: `accuracy = 92%, F1_score = 0.89`

5 Artifacts

- Artifacts are **output files** or **saved results** generated during the run.
- Includes:
 - Trained models (pickle files, .pkl)
 - Plots or graphs
 - Logs
 - Preprocessed data



These can be reused for deployment or further testing.

Example: Trained model file → `student_model.pkl`

⌚ Summary Flow:

text

CopyEdit

Experiment

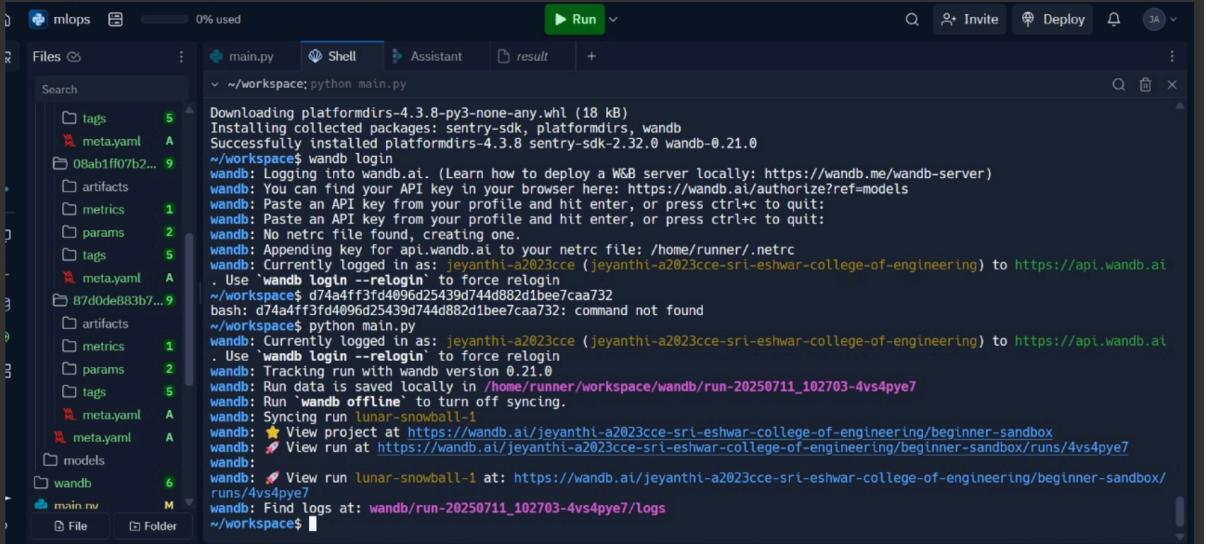
```
  └── Run #1
      ├── Parameters → (e.g. learning_rate=0.01)
      ├── Metrics   → (e.g. accuracy=92%)
      └── Artifacts  → (e.g. saved_model.pkl)
  └── Run #2
      ├── Parameters → ...
      ├── Metrics   → ...
      └── Artifacts  → ...
```



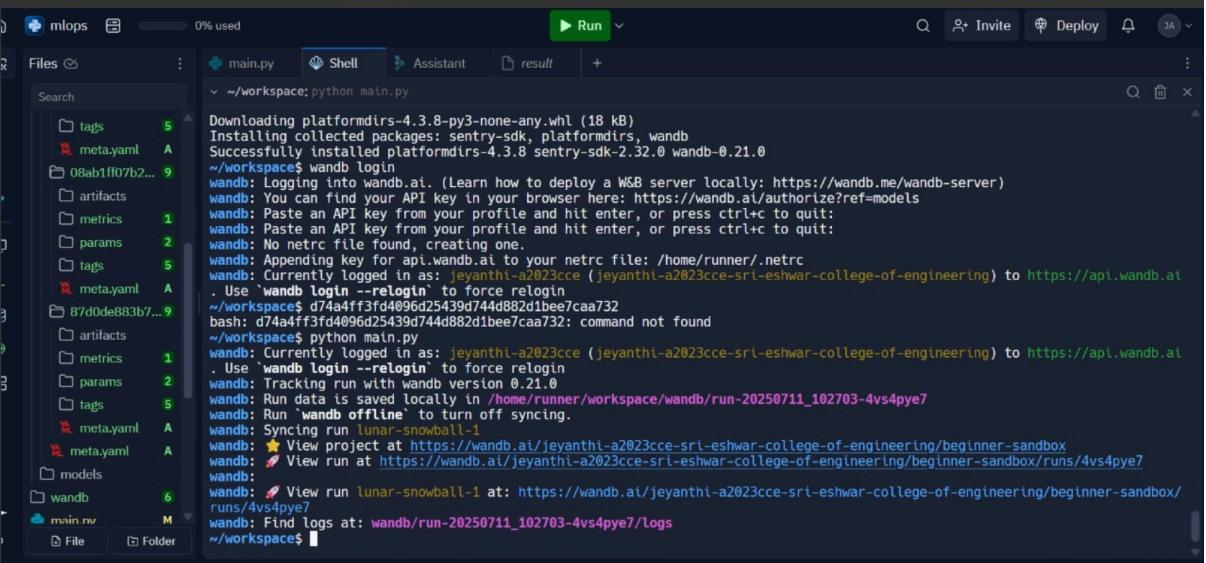
✓ Benefits for / Business Teams:

Feature	Business Value
Experiment tracking	Clear visibility of model improvements
Run comparison	Helps choose the best-performing model
Parameter & metric logging	Scientific, reproducible ML approach
Artifact storage	Easy deployment and audit-ready models

WEIGHTS AND BIASES



```
~/workspace$ wandb login
wandb: Logging into wandb.ai. (Learn how to deploy a W&B server locally: https://wandb.me/wandb-server)
wandb: You can find your API key in your browser here: https://wandb.ai/authorize?ref=models
wandb: Paste an API key from your profile and hit enter, or press ctrl+c to quit:
wandb: Paste an API key from your profile and hit enter, or press ctrl+c to quit:
wandb: No netrc file found, creating one.
wandb: Appending key for api.wandb.ai to your netrc file: /home/runner/.netrc
wandb: Currently logged in as: jeyanthi-a2023cce (jeyanthi-a2023cce-sri-eshwar-college-of-engineering) to https://api.wandb.ai
. Use 'wandb login --relogin' to force relogin
~/workspace$ d74a4ff3fd4096d25439d744d882d1bee7caa732
bash: d74a4ff3fd4096d25439d744d882d1bee7caa732: command not found
~/workspace$ python main.py
wandb: Currently logged in as: jeyanthi-a2023cce (jeyanthi-a2023cce-sri-eshwar-college-of-engineering) to https://api.wandb.ai
. Use 'wandb login --relogin' to force relogin
wandb: Tracking run with wandb version 0.21.0
wandb: Run data is saved locally in /home/runner/workspace/wandb/run-20250711_102703-4vs4pye7
wandb: Run 'wandb offline' to turn off syncing.
wandb: Syncing run lunar-snowball-1
wandb: ★ View project at https://wandb.ai/jeyanthi-a2023cce-sri-eshwar-college-of-engineering/beginner-sandbox
wandb: View run at https://wandb.ai/jeyanthi-a2023cce-sri-eshwar-college-of-engineering/beginner-sandbox/runs/4vs4pye7
wandb:
wandb: 🚫 View run lunar-snowball-1 at: https://wandb.ai/jeyanthi-a2023cce-sri-eshwar-college-of-engineering/beginner-sandbox/runs/4vs4pye7
wandb: Find logs at: wandb/run-20250711_102703-4vs4pye7/logs
~/workspace$
```



```
~/workspace$ wandb login
wandb: Logging into wandb.ai. (Learn how to deploy a W&B server locally: https://wandb.me/wandb-server)
wandb: You can find your API key in your browser here: https://wandb.ai/authorize?ref=models
wandb: Paste an API key from your profile and hit enter, or press ctrl+c to quit:
wandb: Paste an API key from your profile and hit enter, or press ctrl+c to quit:
wandb: No netrc file found, creating one.
wandb: Appending key for api.wandb.ai to your netrc file: /home/runner/.netrc
wandb: Currently logged in as: jeyanthi-a2023cce (jeyanthi-a2023cce-sri-eshwar-college-of-engineering) to https://api.wandb.ai
. Use 'wandb login --relogin' to force relogin
~/workspace$ d74a4ff3fd4096d25439d744d882d1bee7caa732
bash: d74a4ff3fd4096d25439d744d882d1bee7caa732: command not found
~/workspace$ python main.py
wandb: Currently logged in as: jeyanthi-a2023cce (jeyanthi-a2023cce-sri-eshwar-college-of-engineering) to https://api.wandb.ai
. Use 'wandb login --relogin' to force relogin
wandb: Tracking run with wandb version 0.21.0
wandb: Run data is saved locally in /home/runner/workspace/wandb/run-20250711_102703-4vs4pye7
wandb: Run 'wandb offline' to turn off syncing.
wandb: Syncing run lunar-snowball-1
wandb: ★ View project at https://wandb.ai/jeyanthi-a2023cce-sri-eshwar-college-of-engineering/beginner-sandbox
wandb: View run at https://wandb.ai/jeyanthi-a2023cce-sri-eshwar-college-of-engineering/beginner-sandbox/runs/4vs4pye7
wandb:
wandb: 🚫 View run lunar-snowball-1 at: https://wandb.ai/jeyanthi-a2023cce-sri-eshwar-college-of-engineering/beginner-sandbox/runs/4vs4pye7
wandb: Find logs at: wandb/run-20250711_102703-4vs4pye7/logs
~/workspace$
```

The screenshot shows the Wandb interface. On the left, a sidebar lists 'Core' (Registry, Inference), 'Profile' (jeyanthi-a2023cce), and 'Teams'. A 'Pro trial' box indicates 29 days left. On the right, the main area is titled '1. Set up the wandb library'. It shows the command `pip install wandb` and instructions to log in with an API key. A yellow box highlights the API key: `d74a4ff3fd4096d25439d74...`.

This screenshot is identical to the one above, showing the '1. Set up the wandb library' step and the highlighted API key.

The screenshot shows a code editor with a Python file named 'main.py'. The code is as follows:

```
1 import wandb
2 wandb.init(project="beginner-sandbox")
3 a=3
4 b=5
5 result=a+b
6 wandb.log({"result": result})
```

