

Data Versioning (DVC)

→ Why Data Versioning?

- Suppose you built a ML Model
- ↳ Collect data → data-V1.csv
 - ↳ Clean it → data-V1.csv
 - ↳ Add new feature → data-V3.csv (previous code)
 - ↳ Train a model → model-V1.pkl
 - ↳ By new Params → model-V2.pkl
 - ↳ Then you came back next month with a new version of data give best Model

Real world Problem with DVC

- Data Trained a best model, but didn't know which data used
- Your team made don't have access to best data
- You need to track data + model + Params But git only track code.

So why DVC?

DVC integrates with Git

Versioning	✓
Handles large file	✓
Team Sharing	✓
Track Experiment	✗

What is DVC

It is Git but for large Dataset and ML Pipelines

- Version Control datasets and model files (without bloating)
- Track with Version of data for which version of code.
- Share with Team / Cloud easily.

are Simple but yr.

git has changes per version barakar

but DVC uses

How DVC works.

(1) Theory → for every version of git code DVC Maintains Data

Git

DVC

Code
txt
ignore
data

V1

code
txt
ignore
dvc-1

V2

code
txt
ignore
dvc-2

Data

Basically git has large file so, it just store reference of data not Particular DVC

Now DVC Coding.

- Made a file and defined a Simple df. and perfer git.
- Since I have my git repo in Parent Dir MLops, what I will is that I will just track my Day3-DVC folder.
- DVC Install and DVC init →
.dvc, .dvcignore.
- Creating a Space where DVC will store Models/Data.
- As of know I will use my local and make a dir 'S3'
- Used DVC remote add -d myremote S3
- Now we will give folder/file to tracked by DVC
→ we need to detack that file/folder first from git
→ git rm -r --cache file/folder.name
→ git commit -m 'Rem2'

→ now we use (dvc add file/folder) to track

then,

→ in my case I will track Day3DVC/DataDir

→ After doing this we got two file

→ DataDir.dvc → contain id of version

→ .gitignore → contain DataDir

→ we can add more two file
in our git to track changes

→ Now 'dvc commit' and 'dvc push' → This will push our data in our remote space [in my case S3].

↳ Here my S3 folder, contains ^{with} file [for file push 2 file]
↳ 1. file contains 1d
↳ 2. file contains Data.txt

→ Now technically you should git your first version.

→ Now we can make v1 of data.

↳ Change your data. (dvc add my_data.txt)
↳ dvc commit & push
↳ See dvc status
↳ git push origin master

Now our main dvc uses to roll back to version and version

→ List all version → git log --oneline (Copy shield of V1090)

→ Go to that version → git checkout (Should)

or Pull & see Check dvc status → you see Change out

Now Pull that particular version of data → dvc pull

or

git checkout

dvc pull