

* DVC (Data Version Control)

↳ basics

↳ Where Applicable

↳ projects → basics

↳ Sklearn ML

↳ Tensorflow (Deep learning)

↳ Pytorch

↳ This is helpful for Experiment tracking

↳ Data Versioning

↳ model Versioning artifacts.

↳ Automate ML pipelines with DVC

↳ Stage 01 load & save.py

↳ stage 02 split data.py.

↳ Stage 03 train.py

↳ stage 04 evaluate.py

↳ large data files,

ML models,

config files,

scripts & notebooks

→ DVC →

Version Controlled Copies
meta data & info about
changes,
integrates with Git.

↳ it build pipelines.

* DVC (Data Version Control)

↳ basics

↳ Where Applicable

↳ projects → basics

↳ Sklearn ML

↳ Tensorflow (Deep learning)

↳ Pytorch

↳ This is helpful for Experiment tracking

↳ Data Versioning

↳ model Versioning artifacts.

↳ Automate ML pipelines with DVC

↳ Stage 01 load & save.py

↳ stage 02 split data.py.

↳ Stage 03 train.py

↳ stage 04 evaluate.py

↳ large data files,

ML models,

config files,

scripts & notebooks

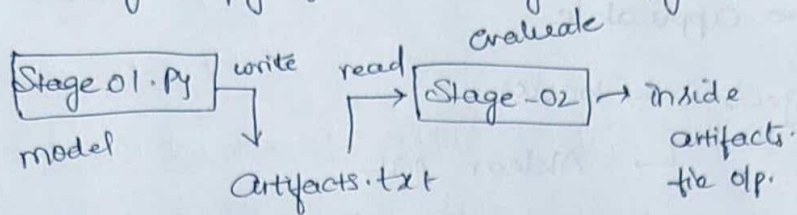
→ DVC →

Version Controlled Copies
meta data & info about
changes,
integrates with Git.

↳ it build pipelines.

→ ~~Install DVC~~

→ Created Stage-01.py file & Stage02.py file.

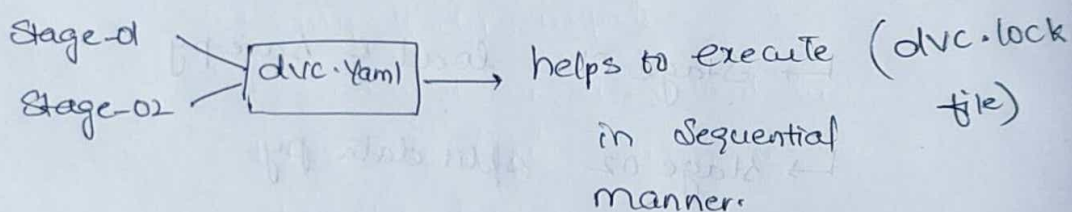


→ in most of ml pipeline we write something & tries to create a file. & then read it.

→ DAG Direct Acyclic graph should follow here we write create & then write. this flow should follow read

→ create .dvc, dvc ignore files (we are initializing them)
dvc init, git init

→ Create dvc.yaml file



→ in dvc.yaml we give

Stages

→ Stage 01

→ cmd (running Stage 01 file)

→ dependencies (files needed to run cmd)

→ O/p of file Stage 01

→ Stage 02

→ cmd

→ dependency

→ O/p's if any.

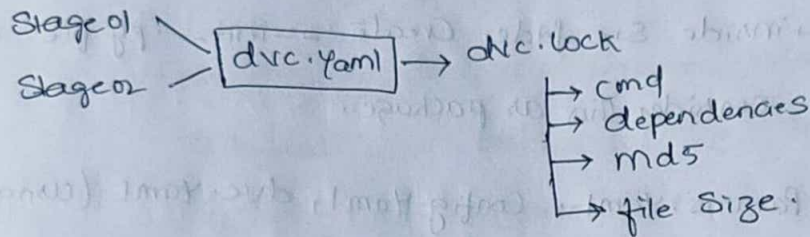
→ ...

at the end this creates dvc.lock file

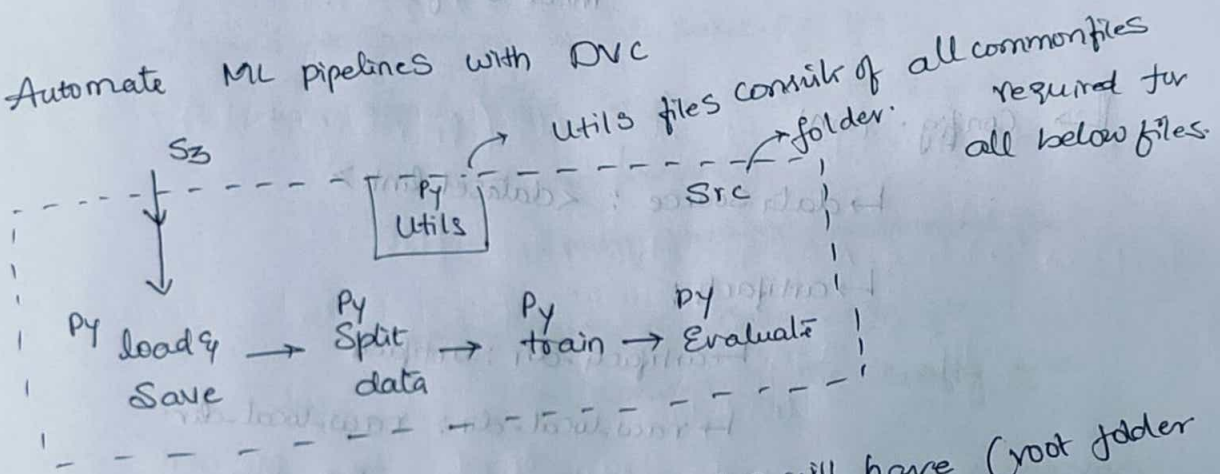
↳ dvc.lock file

↳ it consists of pipeline of stages.

↳ in this along with given variables in dvc.yaml file it has or generate mds (which is a hash value/key) & also size of file it has.



↳ Automate ML pipelines with DVC



↳ Outside of src folder we will have (root folder)

↳ Artifacts folder

↳ it consists raw data, processed data, model, stores metrics,

↳ params.yaml

↳ responsible for having all parameters responsible for training,

↳ dvc.yaml

↳ it consists of all stages steps that occurred, we mention

those here

→ Stage 01. load & Save.py

→ Stage 02 splitdata.py

→ Stage 03 train.py

→ ~~Evaluate~~ Stage 04 Evaluate.py

↳ Config.yaml

↳ After that we Create Requirements.txt file.

↳ Working on DVC pipeline

↳ Create Src folder

↳ In this Create Utils folders.

↳ In Utils Create `--init--.py` file

↳ Inside Src folder Create `--init--.py` file So that it consider this as package.

↳ Create Params.Yaml, Config.Yaml, dvc.Yaml. (Using touch cmd)
inside Config folder

↳ Config.Yaml

↳ data source : <dataset link>

↳ artifact :

↳ artifact-dir : artifacts

↳ raw-local-dir : raw-local-dir

↳ raw-local-file : data.csv

↳ ~~split-data :~~

↳ In Src folder

↳ Create load Save.py file

↳ Inside Utils folder Create utils.py file.

↳ import yaml

↳ import os

↳ `if` path to yaml o/p : dictionary

↳ This takes config.yaml file as `if` & reads it & gives o/p as dictionary form of content in config.Yaml.

→ and then we need to create a Setup.py file.

→ ilp: readme.md file (it reads the file)
& Setup(.....

→ olp: dependency-links.txt

pkg-info

requires.txt

Sources.txt

top-level.txt ...

→ This Setup.py file & Src becomes local packages to

to install local packages, -e. in requirements.txt.

& then run requirements.txt.

→ Now we need to work on load-save.py file.

→ we read ^{yaml} config from all-utils folder from src.

→ ilp: argparse()

Config. Yaml (or) params.yaml file

→ olp: name of given config file

(or) content of the config file.

→ ilp: path to data source in config.yaml

→ olp: it reads the data in path & gives it.

→ Now as we got the data from remote path, we need to save the data that we got from remote path into artifacts directory.

→ Create fn inside this load-save.py file itself.

→ Before create a directory, write a fn in all_utils.py

→ ilp: empty list

→ olp: directories.

→ Now Considering Config.Yaml file
we need to Save dataset in local dir
for that creating path to directory.

→ I/P: artifacts-dir
raw-local-dir
raw-local-file
raw-local-dir-path
raw-local-file-path.

O/P: directory gets Created

& data.csv file gets Created.

→ Till now we have finished stage 01 (loading data)

→ Now adding stage to dvc.Yaml file

I/P: Stages:

load-data:

Cmd: run load-save.py file with Config.Yaml
as I/P

dependencies: all-utils.py

Config.Yaml

load-save.py

O/P: data.csv.

→ commit changes if needed.

→ Now let us go further & add next stage (02) which is splitting of data.

→ Create file in src folder `src/Stage-02-Split-data.py`.

→ i/p: Config path, params path
Split ratio, random-state

O/p:- ① df from data.csv from path

→ fn: Save-local-df → i/p: data, data-path
O/p:- Save file @ datapath

→ datapath → i/p: artifacts_dir, Splitdata_dir, filename

→ Create -directory → i/p: artifacts_dir, split-data-dir

→ Once we do all this O/p:- train.csv

(also) test.csv

→ Now we can add this to dvc.yaml stage.

→ mention deps, params, outs & all in dvc.yaml stage.

→ Save-local-df → i/p: data, datapath
O/p:- datapath.

→ defining directory path for storing split data from config.

→ Now Create stage-03-train.py using touch cmd in src folder and utils folder.

→ In train we need to specify directories & paths along with hyper parameters

→ Create a trained model

→ & Save the model to disk.

→ So for stage 03-train.py

#!/p> go to train model

paths & directories

CSV file

model params

#!/p> trained model in model-dir

(- - - . model file)

→ Now Create another file ~~eva~~ Stage-03 - evaluate.py in src folder

→ i/p :-	trained - model	o/p :- predicted values
	test-data	rmse, mae, r (metrics)
	testy -	Scores dir
	predicted values	Scores filename
	actual values	Scores dir path
	reports-dir	Scores - filepath
	Scores	
	Scores dir	
	Scores - filename	

→ if we want to add a stage where it gives best scores & best model among multiple models.

