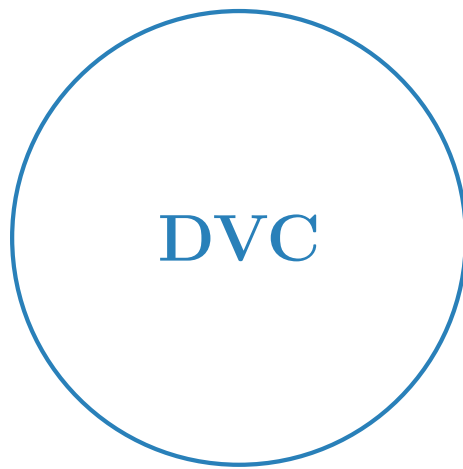# Data Version Control

## DVC Pipeline Configuration Guide

**DVC**

Complete Reference for ML Pipeline Automation

Version Control for Data & Models

**Sujil S**

sujil9480@gmail.com

December 25, 2025

# Contents

# 1 Introduction to DVC

## 1.1 What is DVC?

> **DVC (Data Version Control)** is an open-source version control system for machine learning projects that extends Git's capabilities to handle large datasets, models, and ML workflows.

### 1.1.1 Core Features

- **Large File Management**: Efficiently version datasets >100MB

- **Pipeline Automation**: Reproducible ML workflows

- **Experiment Tracking**: Compare metrics and parameters

- **Model Versioning**: Track model evolution over time

- **Remote Storage**: Works with S3, GCS, Azure, SSH

- **Git Integration**: Seamless integration with Git workflow

## 1.2 Why Use DVC?

### 1.2.1 Traditional Git Limitations

**Problems with Git for ML Projects:**

- Cannot handle files >100MB efficiently

- Not optimized for binary files (models, datasets)

- No built-in experiment tracking

- Lacks pipeline automation

- Difficult to reproduce ML experiments

- Poor performance with large repositories

### 1.2.2 DVC Solutions

**How DVC Solves These Problems:**

- **Efficient Storage**: Links to large files without storing in Git

- **Pipeline Automation**: Define workflows with `dvc.yaml`

- **Experiment Tracking**: Built-in metrics and parameter comparison

- **Reproducibility**: Guaranteed consistent results with `dvc.lock`

- **Collaboration**: Share data seamlessly via remote storage

- **Storage Agnostic**: Compatible with any cloud provider

## 1.3    DVC Architecture

### 1.3.1    Core Components

1. **.dvc files**: Lightweight pointers to actual data

2. **dvc.yaml**: Pipeline definition and stage configuration

3. **dvc.lock**: Lock file with checksums for reproducibility

4. **params.yaml**: Centralized parameter management

5. **metrics/**: JSON/YAML files with evaluation metrics

6. **.dvc/cache**: Local cache for tracked files

7. **Remote Storage**: Cloud storage for team collaboration

### 1.3.2    Workflow Comparison

| Git Workflow | DVC Workflow |
|---|---|
| Tracks source code | Tracks data and models |
| Optimized for text files | Optimized for large binary files |
| Stores in `.git/` directory | Stores in `.dvc/cache` |
| `git commit, push, pull` | `dvc add, push, pull` |
| Remote: GitHub, GitLab | Remote: S3, GCS, Azure |
| Files tracked directly | Files tracked via pointers |

## 1.4    Getting Started

### 1.4.1    Installation

```
# Install DVC with pip
pip install dvc

# Install with specific remote support
pip install dvc[s3]      # AWS S3
pip install dvc[gs]      # Google Cloud Storage
pip install dvc[azure]   # Azure Blob Storage
pip install dvc[ssh]     # SSH remote
pip install dvc[all]     # All remotes

# Verify installation
dvc version
```

### 1.4.2    Project Initialization

```
# Initialize Git repository
git init

# Initialize DVC
dvc init

# Commit DVC configuration
```

```
git add .dvc .gitignore
git commit -m "Initialize DVC"
```

# 2    Understanding dvc.yaml

## 2.1    What is dvc.yaml?

> **dvc.yaml** is the central pipeline configuration file that defines your entire ML workflow. It describes stages, dependencies, commands, outputs, parameters, and metrics in a declarative YAML format.

## 2.2    Basic Structure

### 2.2.1    Minimal Pipeline

```
1  stages:
2    stage_name:
3      cmd: python script.py
4      deps:
5        - input_file.csv
6      outs:
7        - output_file.csv
```

### 2.2.2    Complete Stage Definition

```
1   stages:
2     training:
3       cmd: python train.py --epochs 100
4       wdir: src/              # Working directory
5       deps:
6         - data/processed/train.csv   # Data dependencies
7         - src/train.py          # Code dependencies
8         - src/models/architecture.py  # Module dependencies
9       params:
10        - training.learning_rate     # Parameter references
11        - training.batch_size
12        - model.architecture
13      outs:
14        - models/model.pkl        # Model outputs
15        - models/weights.h5
16      metrics:
17        - metrics/scores.json:     # Metrics (not cached)
18          cache: false
19      plots:
20        - plots/training_curve.csv:   # Plot data
21          cache: false
22          x: epoch
23          y: loss
24      frozen: false            # Allow automatic execution
```

## 2.3 Key Benefits

> **Why Use dvc.yaml?**
>
> - **Reproducibility**: Execute exact same workflow every time
>
> - **Automation**: Run entire pipeline with single command
>
> - **Dependency Tracking**: Only re-run stages with changes
>
> - **Version Control**: Track pipeline evolution in Git
>
> - **Documentation**: Self-documenting workflow
>
> - **Collaboration**: Share workflows with team members

## 2.3 Key Benefits

# 3 Core Components

## 3.1 Stages

**Stages** are the fundamental building blocks representing discrete steps in your ML pipeline.

### 3.1.1 Stage Characteristics

- Each stage has a unique identifier name

- Stages execute in dependency order (DAG)

- Outputs are cached to avoid redundant computation

- Changes in dependencies trigger re-execution

- Multiple stages can run in parallel if independent

### 3.1.2 Simple Stage Example

```
stages:
 preprocess:
   cmd: python src/preprocess.py
   deps:
     - data/raw/dataset.csv
     - src/preprocess.py
   outs:
     - data/processed/train.csv
     - data/processed/test.csv
```

## 3.2 Commands (cmd)

The **cmd** field specifies the shell command to execute.

### 3.2.1 Command Variations

```
# Simple command
stages:
 train:
   cmd: python train.py

# Command with arguments
stages:
 train:
   cmd: python train.py --epochs 100 --lr 0.001

# Multi-line command
stages:
 pipeline:
   cmd: >
     python preprocess.py &&
     python train.py &&
     python evaluate.py

# Command with environment variables
```

```
20  stages:
21   gpu_train:
22     cmd: CUDA_VISIBLE_DEVICES=0 python train.py
23
24  # Command with shell features
25  stages:
26   batch_process:
27     cmd: for file in data/*.csv; do python process.py $file; done
```

## 3.3 Dependencies (deps)

**Dependencies** are files or directories that trigger stage re-execution when modified.

### 3.3.1 Dependency Types

```
1  stages:
2   train:
3     cmd: python train.py
4     deps:
5       # Data dependencies
6       - data/processed/train.csv
7       - data/processed/validation.csv
8
9       # Code dependencies
10      - src/train.py
11      - src/models/neural_network.py
12      - src/utils/data_loader.py
13
14      # Configuration files
15      - configs/model_config.yaml
16
17      # Pre-trained models
18      - models/pretrained/base_model.pkl
19
20      # Directory dependencies
21      - data/images/
```

> **Important**: Always include ALL files that affect stage output. Missing dependencies can lead to inconsistent results and break reproducibility.

## 3.4 Outputs (outs)

**Outputs** are files or directories produced by a stage.

### 3.4.1 Output Configuration

```
1  stages:
2   train:
3     cmd: python train.py
4     outs:
5       # Simple outputs (cached by default)
6       - models/model.pkl
7       - models/weights.h5
```

```
 8
 9     # Outputs with custom settings
10     - models/final_model.pkl:
11       cache: true        # Cache this output
12       persist: false      # Remove when not current
13
14     - logs/training.log:
15       cache: false        # Don't cache logs
16
17     - checkpoints/:
18       cache: true
19       persist: true       # Keep even when not current
```

### 3.4.2   Output Options

| Option | Description |
|---|---|
| `cache` | Cache output (default: `true`) |
| `persist` | Keep in workspace when not current (default: `false`) |
| `checkpoint` | Mark as ML checkpoint for experiments |
| `desc` | Human-readable description |

## 3.5   Parameters (params)

**Parameters** reference values from external files (typically `params.yaml`).

### 3.5.1   Parameter Usage

```
 1  # dvc.yaml
 2  stages:
 3   train:
 4    cmd: python train.py
 5    params:
 6      - training.epochs       # Specific parameter
 7      - training.learning_rate
 8      - training.batch_size
 9      - model            # Entire section
10      - optimizer.type
11      - optimizer.beta1
```

```
 1  # params.yaml
 2  training:
 3   epochs: 100
 4   learning_rate: 0.001
 5   batch_size: 32
 6   early_stopping: true
 7
 8  model:
 9   architecture: resnet50
10   layers: 18
11   dropout: 0.5
12   activation: relu
13
14  optimizer:
15   type: adam
16   beta1: 0.9
```

```
17   beta2: 0.999
```

### 3.5.2   Multiple Parameter Files

```
1  stages:
2   train:
3    cmd: python train.py
4    params:
5      - params.yaml:
6        - training
7        - model
8      - configs/advanced.yaml:
9        - augmentation
10       - preprocessing
```

## 3.6   Metrics

**Metrics** are evaluation results stored in JSON, YAML, CSV, or TSV format.

### 3.6.1   Metric Definition

```
1  stages:
2   evaluate:
3    cmd: python evaluate.py
4    deps:
5      - data/test.csv
6      - models/model.pkl
7    metrics:
8      - metrics/test_scores.json:
9        cache: false       # NEVER cache metrics
10     - metrics/confusion_matrix.json:
11       cache: false
```

### 3.6.2   Metric File Format

```
1  {
2    "accuracy": 0.9542,
3    "precision": 0.9321,
4    "recall": 0.9456,
5    "f1_score": 0.9388,
6    "auc_roc": 0.9782,
7    "confusion_matrix": {
8      "true_positive": 850,
9      "false_positive": 45,
10     "true_negative": 920,
11     "false_negative": 35
12   }
13  }
```

### 3.6.3 Viewing Metrics

```
# Show current metrics
dvc metrics show

# Compare with another branch
dvc metrics diff main

# Compare with previous commit
dvc metrics diff HEAD~1

# Show all experiments
dvc metrics show -R
```

## 3.7 Plots

**Plots** are data files for visualization (CSV, JSON, YAML).

### 3.7.1 Plot Configuration

```
1  stages:
2   evaluate:
3    cmd: python evaluate.py
4    plots:
5      # Simple plot
6      - plots/training_history.csv:
7        cache: false
8
9      # Customized plot
10     - plots/roc_curve.json:
11       cache: false
12       x: fpr
13       y: tpr
14       title: "ROC Curve"
15       x_label: "False Positive Rate"
16       y_label: "True Positive Rate"
17
18     # Multi-series plot
19     - plots/metrics.csv:
20       x: epoch
21       y:
22        plots/metrics.csv: [loss, val_loss]
23       title: "Training vs Validation Loss"
```

### 3.7.2 Plot Data Format

```
1  # plots/training_history.csv
2  epoch,loss,accuracy,val_loss,val_accuracy
3  1,0.693,0.501,0.685,0.515
4  2,0.612,0.653,0.598,0.672
5  3,0.523,0.745,0.501,0.758
6  4,0.445,0.812,0.432,0.825
7  5,0.389,0.856,0.398,0.861
```

# 4 Advanced Pipeline Features

## 4.1 Working Directory (wdir)

Execute stages in a specific working directory.

```
 1  stages:
 2   notebook_analysis:
 3    cmd: jupyter nbconvert --execute analysis.ipynb
 4    wdir: notebooks/
 5    deps:
 6     - notebooks/analysis.ipynb
 7     - data/input.csv
 8    outs:
 9     - notebooks/output.html
10     - notebooks/results.csv
```

## 4.2 Frozen Stages

Prevent automatic execution of specific stages.

```
 1  stages:
 2   expensive_preprocessing:
 3    cmd: python expensive_process.py
 4    frozen: true          # Skip in dvc repro
 5    deps:
 6     - data/raw/large_dataset.csv
 7    outs:
 8     - data/processed/features.pkl
```

> **Use Cases for Frozen Stages:**
>
> - Computationally expensive operations
>
> - Stages requiring manual intervention
>
> - Optional pipeline branches
>
> - Data collection/download stages

## 4.3 Always Changed Stages

Force re-execution regardless of dependencies.

```
 1  stages:
 2   fetch_data:
 3    cmd: python fetch_latest_data.py
 4    always_changed: true     # Always runs
 5    outs:
 6     - data/latest/dataset.csv
 7     - data/latest/metadata.json
```

## 4.4 Foreach Loops

Execute the same stage multiple times with different parameters.

### 4.4.1   Simple Foreach

```
1  stages:
2   train_models:
3    foreach:
4     - logistic_regression
5     - random_forest
6     - gradient_boosting
7     - neural_network
8     - svm
9    do:
10    cmd: python train.py --model ${item}
11    deps:
12     - data/train.csv
13     - src/train.py
14    params:
15     - models.${item}
16    outs:
17     - models/${item}/model.pkl
18    metrics:
19     - metrics/${item}_scores.json:
20       cache: false
```

### 4.4.2   Dictionary Foreach

```
1  stages:
2   process_datasets:
3    foreach:
4     train: data/raw/train.csv
5     validation: data/raw/val.csv
6     test: data/raw/test.csv
7    do:
8     cmd: python process.py --input ${item} --output data/processed/${key}.csv
9     deps:
10     - ${item}
11     - src/process.py
12    outs:
13     - data/processed/${key}.csv
```

### 4.4.3   Matrix Foreach (Hyperparameter Grid)

```
1  stages:
2   grid_search:
3    foreach:
4     lr_0001_bs16:
5      learning_rate: 0.001
6      batch_size: 16
7     lr_0001_bs32:
8      learning_rate: 0.001
9      batch_size: 32
10    lr_001_bs16:
11     learning_rate: 0.01
12     batch_size: 16
13    lr_001_bs32:
14     learning_rate: 0.01
```

```
15        batch_size: 32
16    do:
17     cmd: >
18      python train.py
19      --lr ${item.learning_rate}
20      --batch ${item.batch_size}
21     deps:
22      - data/train.csv
23     outs:
24      - models/${key}_model.pkl
25     metrics:
26      - metrics/${key}_scores.json:
27        cache: false
```

## 4.5   Variables and Templating

Define reusable variables for cleaner configuration.

```
1  vars:
2   - base_dir: .
3   - data_dir: data
4   - models_dir: models
5   - metrics_dir: metrics
6   - src_dir: src
7   - batch_size: 32
8   - learning_rate: 0.001
9   - random_seed: 42
10
11 stages:
12  preprocess:
13    cmd: python ${src_dir}/preprocess.py --seed ${random_seed}
14    deps:
15     - ${data_dir}/raw/dataset.csv
16     - ${src_dir}/preprocess.py
17    outs:
18     - ${data_dir}/processed/train.csv
19     - ${data_dir}/processed/test.csv
20
21  train:
22    cmd: >
23     python ${src_dir}/train.py
24     --batch-size ${batch_size}
25     --lr ${learning_rate}
26     --seed ${random_seed}
27    deps:
28     - ${data_dir}/processed/train.csv
29     - ${src_dir}/train.py
30    outs:
31     - ${models_dir}/model.pkl
32    metrics:
33     - ${metrics_dir}/scores.json:
34       cache: false
```

## 4.6   Stage Dependencies

Explicitly define execution order (rarely needed).

```
1   stages:
2     stage_a:
3       cmd: python a.py
4       outs:
5         - output_a.csv
6
7     stage_b:
8       cmd: python b.py
9       deps:
10        - output_a.csv      # Implicit dependency on stage_a
11      outs:
12        - output_b.csv
13
14    stage_c:
15      cmd: python c.py
16      deps:
17        - output_b.csv      # Implicit dependency on stage_b
18      outs:
19        - output_c.csv
```

**Note**: DVC automatically determines execution order from file dependencies. Explicit stage dependencies are rarely needed.

# 5 Complete ML Pipeline Example

## 5.1 Project Structure

```
ml_project/
|-- dvc.yaml
|-- params.yaml
|-- .dvc/
|   |-- config
|   '-- cache/
|-- data/
|   |-- raw/
|   |-- processed/
|   '-- features/
|-- models/
|-- metrics/
|-- plots/
|-- src/
|   |-- data_collection.py
|   |-- preprocessing.py
|   |-- feature_engineering.py
|   |-- train.py
|   |-- evaluate.py
|   '-- models/
|       '-- architectures.py
|-- notebooks/
|-- requirements.txt
'-- README.md
```

## 5.2 Complete dvc.yaml

```yaml
1  vars:
2    - data_dir: data
3    - models_dir: models
4    - metrics_dir: metrics
5    - plots_dir: plots
6    - src_dir: src
7
8  stages:
9    # =======================================
10   # Stage 1: Data Collection
11   # =======================================
12   collect_data:
13     cmd: python ${src_dir}/data_collection.py
14     deps:
15       - ${src_dir}/data_collection.py
16     params:
17       - data_collection.source_url
18       - data_collection.date_range
19       - data_collection.api_version
20     outs:
21       - ${data_dir}/raw/dataset.csv
22       - ${data_dir}/raw/metadata.json
```

```
23
24    # ==========================================
25    # Stage 2: Data Preprocessing
26    # ==========================================
27    preprocess:
28      cmd: python ${src_dir}/preprocessing.py
29      deps:
30        - ${data_dir}/raw/dataset.csv
31        - ${src_dir}/preprocessing.py
32      params:
33        - preprocessing.train_split
34        - preprocessing.test_split
35        - preprocessing.validation_split
36        - preprocessing.random_seed
37        - preprocessing.handle_missing
38        - preprocessing.remove_outliers
39      outs:
40        - ${data_dir}/processed/train.csv
41        - ${data_dir}/processed/test.csv
42        - ${data_dir}/processed/validation.csv
43        - ${data_dir}/processed/stats.json
44
45    # ==========================================
46    # Stage 3: Feature Engineering
47    # ==========================================
48    feature_engineering:
49      cmd: python ${src_dir}/feature_engineering.py
50      deps:
51        - ${data_dir}/processed/train.csv
52        - ${data_dir}/processed/test.csv
53        - ${data_dir}/processed/validation.csv
54        - ${src_dir}/feature_engineering.py
55      params:
56        - features.numerical_features
57        - features.categorical_features
58        - features.text_features
59        - features.scaling_method
60        - features.encoding_method
61      outs:
62        - ${data_dir}/features/X_train.npy
63        - ${data_dir}/features/X_test.npy
64        - ${data_dir}/features/X_val.npy
65        - ${data_dir}/features/y_train.npy
66        - ${data_dir}/features/y_test.npy
67        - ${data_dir}/features/y_val.npy
68        - ${models_dir}/preprocessors/scaler.pkl
69        - ${models_dir}/preprocessors/encoder.pkl
70
71    # ==========================================
72    # Stage 4: Model Training
73    # ==========================================
74    train:
75      cmd: python ${src_dir}/train.py
76      deps:
77        - ${data_dir}/features/X_train.npy
78        - ${data_dir}/features/X_val.npy
79        - ${data_dir}/features/y_train.npy
80        - ${data_dir}/features/y_val.npy
```

```yaml
81      - ${src_dir}/train.py
82      - ${src_dir}/models/architectures.py
83    params:
84      - training.model_type
85      - training.epochs
86      - training.batch_size
87      - training.learning_rate
88      - training.optimizer
89      - training.loss_function
90      - training.early_stopping
91      - training.random_seed
92    outs:
93      - ${models_dir}/trained/model.pkl
94      - ${models_dir}/trained/weights.h5
95      - ${models_dir}/trained/history.json
96    plots:
97      - ${plots_dir}/training_loss.csv:
98        cache: false
99        x: epoch
100       y: [loss, val_loss]
101       title: "Training vs Validation Loss"
102     - ${plots_dir}/training_accuracy.csv:
103       cache: false
104       x: epoch
105       y: [accuracy, val_accuracy]
106
107   # =========================================
108   # Stage 5: Model Evaluation
109   # =========================================
110   evaluate:
111    cmd: python ${src_dir}/evaluate.py
112    deps:
113      - ${data_dir}/features/X_test.npy
114      - ${data_dir}/features/y_test.npy
115      - ${models_dir}/trained/model.pkl
116      - ${models_dir}/preprocessors/scaler.pkl
117      - ${src_dir}/evaluate.py
118    params:
119      - evaluation.metrics
120      - evaluation.threshold
121    outs:
122      - predictions/test_predictions.csv
123    metrics:
124      - ${metrics_dir}/test_scores.json:
125        cache: false
126    plots:
127      - ${plots_dir}/confusion_matrix.csv:
128        cache: false
129        template: confusion
130      - ${plots_dir}/roc_curve.json:
131        cache: false
132        x: fpr
133        y: tpr
134        title: "ROC Curve"
135      - ${plots_dir}/precision_recall.json:
136        cache: false
137        x: recall
138        y: precision
```

## 5.3 Complete params.yaml

```yaml
 1  # ========================================
 2  # Data Collection Parameters
 3  # ========================================
 4  data_collection:
 5    source_url: "https://api.example.com/v2/dataset"
 6    date_range:
 7      start: "2024-01-01"
 8      end: "2024-12-31"
 9    api_version: "v2.0"
10    timeout: 30
11
12  # ========================================
13  # Preprocessing Parameters
14  # ========================================
15  preprocessing:
16    train_split: 0.70
17    test_split: 0.20
18    validation_split: 0.10
19    random_seed: 42
20    handle_missing: "median"
21    remove_outliers: true
22    outlier_method: "iqr"
23    outlier_threshold: 1.5
24
25  # ========================================
26  # Feature Engineering
27  # ========================================
28  features:
29    numerical_features:
30      - age
31      - income
32      - credit_score
33      - balance
34
35    categorical_features:
36      - gender
37      - occupation
38      - education
39
40    text_features:
41      - description
42
43    scaling_method: "standard"
44    encoding_method: "onehot"
45
46  # ========================================
47  # Training Parameters
48  # ========================================
49  training:
50    model_type: "random_forest"
51    epochs: 100
52    batch_size: 32
53    learning_rate: 0.001
54    optimizer: "adam"
55    loss_function: "binary_crossentropy"
56    early_stopping:
```

```
57      enabled: true
58      patience: 10
59      min_delta: 0.001
60    random_seed: 42
61
62    # =========================================
63    # Evaluation Parameters
64    # =========================================
65    evaluation:
66     metrics:
67       - accuracy
68       - precision
69       - recall
70       - f1_score
71       - roc_auc
72      threshold: 0.5
```

# 6 Essential DVC Commands

## 6.1 Initialization and Setup

```
# Initialize DVC in project
dvc init

# Configure S3 remote
dvc remote add -d storage s3://mybucket/dvcstore

# Configure GCS remote
dvc remote add -d storage gs://mybucket/dvcstore

# List remotes
dvc remote list

# Modify remote URL
dvc remote modify storage url s3://newbucket/path
```

## 6.2 Pipeline Operations

```
# Run entire pipeline
dvc repro

# Run specific stage
dvc repro train

# Force re-run all stages
dvc repro -f

# Force re-run specific stage
dvc repro -f preprocess

# Show pipeline DAG
dvc dag

# Check pipeline status
dvc status

# Detailed status
dvc status -v
```

## 6.3 Data Management

```
# Track large file
dvc add data/dataset.csv

# Track directory
dvc add data/images/
```

```
# Push data to remote
dvc push

# Pull data from remote
dvc pull

# Fetch to cache only
dvc fetch

# Update workspace from cache
dvc checkout
```

## 6.4   Metrics and Plots

```
# Show metrics
dvc metrics show

# Compare metrics with branch
dvc metrics diff main

# Show plots
dvc plots show

# Compare plots
dvc plots diff main experiment

# Generate HTML report
dvc plots show --html
```

## 6.5   Experiments

```
# Run experiment
dvc exp run

# Run with parameter override
dvc exp run -S train.lr=0.01

# Queue experiments
dvc exp run --queue -S train.lr=0.001
dvc exp run --queue -S train.lr=0.01
dvc queue start

# List experiments
dvc exp show

# Compare experiments
dvc exp diff
```

```
# Apply experiment
dvc exp apply exp-12345
```

# 7 Best Practices

## 7.1 Pipeline Design

**Good Practices:**

- Create stages for logical workflow steps

- Keep stages focused (single responsibility)

- Use descriptive stage names (`preprocess_data`, not `stage1`)

- Track all dependencies accurately

- Organize outputs in structured directories

- Use `params.yaml` for all configurable values

**Avoid:**

- Overly granular stages (creates overhead)

- Combining unrelated operations in one stage

- Missing dependencies (breaks reproducibility)

- Hardcoding parameters in commands

- Caching log files and metrics

## 7.2 Version Control Integration

### 7.2.1 Track with Git

- `dvc.yaml` - Pipeline definition

- `dvc.lock` - Lock file

- `params.yaml` - Parameters

- `*.dvc` - File pointers

- `.dvc/config` - Configuration

- `src/` - Source code

### 7.2.2 Track with DVC

- `data/` - Datasets

- `models/` - Trained models

- Large files (>1MB)

## 7.3   Parameter Management

```
1   # Good: Hierarchical structure
2   preprocessing:
3     train_split: 0.7
4     random_seed: 42
5
6   training:
7     epochs: 100
8     learning_rate: 0.001
9
10  # Bad: Flat structure
11  train_split: 0.7
12  epochs: 100
13  learning_rate: 0.001
```

# 8 Troubleshooting

## 8.1 Pipeline Won't Run

**Problem**: `dvc repro` doesn't execute stages

> **Possible Causes**:

- Pipeline already up-to-date (no changes detected)

- Syntax errors in `dvc.yaml`

- Missing dependencies

- Circular dependencies

- Frozen stages

```
# Check pipeline status
dvc status
dvc status -v

# Validate pipeline structure
dvc dag

# Force re-run entire pipeline
dvc repro -f

# Force re-run specific stage
dvc repro -f stage_name

# Check for syntax errors
cat dvc.yaml | python -m yaml

# Dry run (show what would execute)
dvc repro --dry
```

## 8.2 Stage Always Re-runs

**Problem**: Stage executes every time despite no changes

> **Possible Causes**:

- Missing dependencies in `deps` field

- Files being modified by external processes

- Non-deterministic code (random operations without seed)

- Timestamp issues

- Command generates different output each time

> **Solutions**:

```
# Check what DVC detects as changed
dvc status -v

# Verify file checksums
dvc status --show-json

# Set random seeds in your code
# Python example:
import random
import numpy as np
random.seed(42)
np.random.seed(42)
```

**Ensure Reproducibility**:

- List ALL dependencies (data, code, configs)

- Set random seeds for all random operations

- Avoid timestamp-based operations

- Use deterministic algorithms

## 8.3  Missing Data Files

**Problem**: Data files not found after cloning repository

```
# Pull all data from remote
dvc pull

# Pull specific file
dvc pull data/dataset.csv.dvc

# Check remote configuration
dvc remote list
dvc remote list --show-origin

# Test connectivity to remote
dvc status --cloud

# Check if files are in cache
ls -la .dvc/cache/

# Fetch to cache without checking out
dvc fetch
dvc checkout
```

## 8.4  Push/Pull Failures

**Problem**: Cannot push or pull data to/from remote

    **Common Issues**:

- Incorrect remote URL

- Missing credentials

- Permission issues

- Network connectivity

- Insufficient storage quota

```
# Check remote configuration
dvc remote list
dvc config remote.storage.url

# Verbose output for debugging
dvc push -v
dvc pull -v

# Check cloud status
dvc status --cloud

# For AWS S3 - verify credentials
aws configure list
aws s3 ls s3://your-bucket/

# For Google Cloud - verify authentication
gcloud auth list
gsutil ls gs://your-bucket/

# Test with small file first
dvc add test.txt
dvc push test.txt.dvc
```

## 8.5   Cache Corruption

**Problem**: DVC cache becomes corrupted or inconsistent

```
# Check cache integrity
ls -la .dvc/cache/

# Remove corrupted cache
rm -rf .dvc/cache

# Re-download from remote
dvc fetch
dvc checkout

# Or re-run pipeline
dvc repro -f

# Verify data integrity
dvc status
```

## 8.6   Disk Space Issues

**Problem**: DVC cache consuming too much disk space

```
# Check cache size
du -sh .dvc/cache
dvc cache dir

# Show what would be removed (dry run)
dvc gc --dry

# Remove files not in current workspace
dvc gc --workspace

# Keep workspace, all branches, and all tags
dvc gc --workspace --all-branches --all-tags

# Remove all except workspace
dvc gc --workspace --force

# Clean cloud cache
dvc gc --cloud

# Check after cleanup
du -sh .dvc/cache
```

> **Warning**: Be careful with `dvc gc`. Always use `--dry` first to see what would be removed. Keep backups of important data.

## 8.7   Lock File Conflicts

**Problem**: Git merge conflict in `dvc.lock`

**Solution 1: Accept One Version**

```
# Accept your version
git checkout --ours dvc.lock

# Or accept their version
git checkout --theirs dvc.lock

# Regenerate lock file
dvc repro

# Commit resolution
git add dvc.lock
git commit -m "Resolve dvc.lock conflict"
```

**Solution 2: Regenerate**

```
# Remove conflicted lock file
rm dvc.lock

# Regenerate from pipeline
dvc repro

# Commit new lock file
git add dvc.lock
git commit -m "Regenerate dvc.lock"
```

## 8.8   Permission Issues

**Problem**: Permission denied errors

```
# Fix local cache permissions
chmod -R u+w .dvc/cache

# Fix workspace file permissions
chmod -R u+w data/ models/

# For remote storage (S3)
# Check IAM permissions in AWS Console

# For SSH remotes
ssh user@host "chmod -R u+w /path/to/dvc/storage"

# Check file ownership
ls -la .dvc/cache
ls -la data/
```

## 8.9   Performance Issues

### 8.9.1   Slow Pipeline Execution

**Optimizations**:

- Enable parallel stage execution (if independent)

- Use faster remote storage region

- Optimize data I/O operations in code

- Use data sampling for development

- Profile code for bottlenecks

- Consider stage granularity

- Use SSD for cache directory

### 8.9.2   Slow Push/Pull Operations

```
# Use faster remote region
dvc remote modify storage region us-west-2

# Enable compression
dvc remote modify storage --local compression gzip

# Increase transfer workers
dvc remote modify storage --local upload_max_workers 20
dvc remote modify storage --local download_max_workers 20

# Use jobs parameter
dvc push --jobs 8
dvc pull --jobs 8
```

## 8.10   Import/Module Errors

**Problem**: Python modules not found during stage execution

```
# Verify Python environment
which python
python --version

# Check installed packages
pip list

# Install requirements
pip install -r requirements.txt

# Check PYTHONPATH
echo $PYTHONPATH

# Add to PYTHONPATH in stage command
stages:
  train:
    cmd: PYTHONPATH=src:$PYTHONPATH python src/train.py
```

# 9 Advanced Topics

## 9.1 Remote Storage Configuration

### 9.1.1 AWS S3

```
# Add S3 remote
dvc remote add -d storage s3://mybucket/dvcstore

# Configure region
dvc remote modify storage region us-west-2

# Use AWS profile
dvc remote modify storage profile myprofile

# Use specific credentials (not recommended)
dvc remote modify storage access_key_id YOUR_KEY
dvc remote modify storage secret_access_key YOUR_SECRET

# Enable server-side encryption
dvc remote modify storage sse AES256

# Use custom endpoint (MinIO, etc.)
dvc remote modify storage endpointurl http://localhost:9000

# Enable versioning
dvc remote modify storage version_aware true
```

### 9.1.2 Google Cloud Storage

```
# Add GCS remote
dvc remote add -d storage gs://mybucket/dvcstore

# Configure project
dvc remote modify storage projectname myproject

# Use service account
dvc remote modify storage credentialpath \
    /path/to/credentials.json

# Set default credentials
export GOOGLE_APPLICATION_CREDENTIALS=/path/to/key.json
```

### 9.1.3 Azure Blob Storage

```
# Add Azure remote
dvc remote add -d storage azure://mycontainer/path

# Use connection string
```

```
dvc remote modify storage connection_string \
    "DefaultEndpointsProtocol=https;..."

# Or use account name and key
dvc remote modify storage account_name myaccount
dvc remote modify storage account_key mykey

# Use SAS token
dvc remote modify storage sas_token "?sv=2019..."
```

### 9.1.4   SSH Remote

```
# Add SSH remote
dvc remote add -d storage \
    ssh://user@example.com/path/to/dvc

# Use SSH key
dvc remote modify storage keyfile ~/.ssh/id_rsa

# Use password (not recommended)
dvc remote modify storage password mypassword

# Specify port
dvc remote modify storage port 2222

# Use SSH config
dvc remote modify storage ssh_config ~/.ssh/config
```

## 9.2   CI/CD Integration

### 9.2.1   GitHub Actions

```yaml
1  # .github/workflows/dvc-pipeline.yml
2  name: DVC Pipeline
3
4  on:
5    push:
6      branches: [main]
7    pull_request:
8      branches: [main]
9
10 jobs:
11   run-pipeline:
12     runs-on: ubuntu-latest
13
14     steps:
15       - name: Checkout code
16         uses: actions/checkout@v3
17
18       - name: Setup Python
19         uses: actions/setup-python@v4
20         with:
```

```
21      python-version: '3.10'
22
23    - name: Install dependencies
24      run: |
25       pip install dvc[s3]
26       pip install -r requirements.txt
27
28    - name: Configure AWS credentials
29      uses: aws-actions/configure-aws-credentials@v2
30      with:
31       aws-access-key-id: ${{ secrets.AWS_ACCESS_KEY_ID }}
32       aws-secret-access-key: ${{ secrets.AWS_SECRET_ACCESS_KEY }}
33       aws-region: us-west-2
34
35    - name: Pull DVC data
36      run: dvc pull
37
38    - name: Run pipeline
39      run: dvc repro
40
41    - name: Show metrics
42      run: dvc metrics show
43
44    - name: Push results
45      if: github.ref == 'refs/heads/main'
46      run: dvc push
47
48    - name: Create CML report
49      run: |
50       echo "## Metrics" >> report.md
51       dvc metrics show --md >> report.md
52       cml comment create report.md
```

## 9.2.2   GitLab CI

```
1  # .gitlab-ci.yml
2  image: python:3.10
3
4  stages:
5   - setup
6   - pipeline
7   - report
8
9  variables:
10  DVC_REMOTE: s3://mybucket/dvcstore
11
12 before_script:
13  - pip install dvc[s3]
14  - pip install -r requirements.txt
15
16 setup_dvc:
17  stage: setup
18  script:
19   - dvc remote modify storage access_key_id $AWS_ACCESS_KEY_ID
20   - dvc remote modify storage secret_access_key $AWS_SECRET_ACCESS_KEY
21   - dvc pull
22  cache:
```

```
23    paths:
24      - .dvc/cache
25
26  run_pipeline:
27    stage: pipeline
28    script:
29      - dvc repro
30      - dvc push
31    artifacts:
32      paths:
33        - metrics/
34        - plots/
35      expire_in: 1 week
36
37  generate_report:
38    stage: report
39    script:
40      - dvc metrics show
41      - dvc plots show
42    artifacts:
43      paths:
44        - metrics/
45        - plots/
```

## 9.3   Custom Metrics and Plots

### 9.3.1   Custom Metrics Structure

```
1  {
2    "train": {
3      "accuracy": 0.9542,
4      "loss": 0.1234,
5      "precision": 0.9321,
6      "recall": 0.9456
7    },
8    "validation": {
9      "accuracy": 0.9312,
10     "loss": 0.1567,
11     "precision": 0.9123,
12     "recall": 0.9287
13   },
14   "test": {
15     "accuracy": 0.9278,
16     "loss": 0.1689,
17     "precision": 0.9056,
18     "recall": 0.9201
19   },
20   "model_info": {
21     "parameters": 1250000,
22     "training_time": 3600,
23     "epochs_trained": 50
24   }
25 }
```

### 9.3.2   Custom Plot Templates

```
 1  stages:
 2   evaluate:
 3    plots:
 4      - plots/custom_metrics.csv:
 5        template: linear
 6        x: iteration
 7        y:
 8         plots/custom_metrics.csv: [metric_a, metric_b]
 9        title: "Custom Metrics Comparison"
10        x_label: "Training Iteration"
11        y_label: "Metric Value"
12
13      - plots/scatter_plot.csv:
14        template: scatter
15        x: predicted
16        y: actual
17        title: "Predicted vs Actual"
18
19      - plots/smooth_curve.csv:
20        template: smooth
21        x: epoch
22        y: loss
23        title: "Smoothed Loss Curve"
```

## 9.4   Pipeline Modularization

### 9.4.1   Splitting Large Pipelines

```
project/
|-- dvc.yaml
|-- pipelines/
|   |-- data_pipeline.yaml
|   |-- training_pipeline.yaml
|   '-- evaluation_pipeline.yaml
|-- params/
|   |-- data_params.yaml
|   |-- training_params.yaml
|   '-- evaluation_params.yaml
```

```
 1  # pipelines/data_pipeline.yaml
 2  stages:
 3   collect:
 4    cmd: python src/collect.py
 5    deps:
 6     - src/collect.py
 7    outs:
 8     - data/raw/dataset.csv
 9
10   preprocess:
11    cmd: python src/preprocess.py
12    deps:
13     - data/raw/dataset.csv
14     - src/preprocess.py
15    outs:
```

```
16      - data/processed/train.csv
17      - data/processed/test.csv
```

```
1   # dvc.yaml (main pipeline)
2   stages:
3     data_pipeline:
4       cmd: dvc repro pipelines/data_pipeline.yaml
5
6     training_pipeline:
7       cmd: dvc repro pipelines/training_pipeline.yaml
8       deps:
9         - data/processed/train.csv
10
11    evaluation_pipeline:
12      cmd: dvc repro pipelines/evaluation_pipeline.yaml
13      deps:
14        - models/trained/model.pkl
```

## 9.5  Experiment Management

### 9.5.1  Running Multiple Experiments

```
# Queue multiple experiments
dvc exp run --queue -S train.lr=0.001 -S train.epochs=50
dvc exp run --queue -S train.lr=0.01 -S train.epochs=50
dvc exp run --queue -S train.lr=0.1 -S train.epochs=50
dvc exp run --queue -S train.lr=0.001 -S train.epochs=100


# Start queue execution
dvc queue start


# Run experiments in parallel (4 workers)
dvc queue start --jobs 4


# Check queue status
dvc queue status


# Stop queue
dvc queue stop


# Remove experiments from queue
dvc queue remove exp-12345
```

### 9.5.2  Comparing Experiments

```
# List all experiments
dvc exp show


# Show only changed parameters and metrics
dvc exp show --only-changed
```

```
# Show experiments in table format
dvc exp show --no-pager

# Compare specific experiments
dvc exp diff exp-abc123 exp-def456

# Compare with specific branch
dvc exp diff main

# Show differences in specific metrics
dvc exp show --include-params train.lr \
    --include-metrics accuracy,f1_score
```

### 9.5.3   Managing Experiments

```
# Apply experiment to workspace
dvc exp apply exp-abc123

# Create branch from experiment
dvc exp branch exp-abc123 feature/best-model

# Push experiment to remote
dvc exp push origin exp-abc123

# Pull experiments from remote
dvc exp pull origin

# Remove experiment
dvc exp remove exp-abc123

# Remove all experiments
dvc exp remove --all

# Remove experiments in queue
dvc exp remove --queue
```

# 10 Quick Reference Guide

## 10.1 dvc.yaml Components Summary

| Component | Description |
|---|---|
| `stages` | Define individual pipeline stages |
| `cmd` | Shell command to execute |
| `wdir` | Working directory for execution |
| `deps` | File/directory dependencies |
| `outs` | Output files/directories |
| `params` | Parameter references from `params.yaml` |
| `metrics` | Evaluation metric files (JSON/YAML/CSV) |
| `plots` | Plot data files for visualization |
| `frozen` | Prevent automatic execution |
| `always_changed` | Force re-execution every time |
| `foreach` | Loop over items |
| `vars` | Define reusable variables |

## 10.2 Command Cheat Sheet

**Initialization & Setup**

```
dvc init                    # Initialize DVC
dvc remote add -d name url  # Add default remote
dvc remote list             # List remotes
```

**Pipeline Operations**

```
dvc repro                   # Run pipeline
dvc repro -f                # Force re-run
dvc repro stage_name        # Run specific stage
dvc dag                     # Show pipeline DAG
dvc status                  # Check status
```

**Data Management**

```
dvc add file                # Track file
dvc push                    # Push to remote
dvc pull                    # Pull from remote
dvc fetch                   # Fetch to cache
dvc checkout                # Update workspace
```

**Metrics & Plots**

```
dvc metrics show            # Show metrics
dvc metrics diff            # Compare metrics
dvc plots show              # Show plots
dvc plots diff              # Compare plots
```

> **Experiments**
> ```
> dvc exp run                # Run experiment
> dvc exp show               # List experiments
> dvc exp diff               # Compare experiments
> dvc exp apply exp-id       # Apply experiment
> dvc exp branch exp-id name # Create branch
> ```

## 10.3   Common Workflows

### 10.3.1   Start New DVC Project

```
# Initialize repositories
git init
dvc init

# Configure remote storage
dvc remote add -d storage s3://mybucket/dvcstore

# Commit configuration
git add .dvc .gitignore
git commit -m "Initialize DVC"

# Track large dataset
dvc add data/dataset.csv
git add data/dataset.csv.dvc data/.gitignore
git commit -m "Track dataset"

# Push data to remote
dvc push
```

### 10.3.2   Clone DVC Project

```
# Clone repository
git clone https://github.com/user/repo.git
cd repo

# Pull data from remote
dvc pull

# Verify pipeline
dvc dag
dvc status

# Run pipeline
dvc repro
```

### 10.3.3   Run Experiment with Different Parameters

```
# Create feature branch
git checkout -b experiment/new-model

# Modify parameters
vim params.yaml

# Run pipeline
dvc repro

# Compare with main branch
dvc metrics diff main
dvc plots diff main

# If successful, commit and merge
git add params.yaml dvc.lock
git commit -m "Experiment: improved model"
dvc push
git push origin experiment/new-model
```

### 10.3.4   Hyperparameter Tuning

```
# Queue multiple experiments
for lr in 0.001 0.01 0.1; do
  for bs in 16 32 64; do
    dvc exp run --queue \
      -S train.learning_rate=$lr \
      -S train.batch_size=$bs \
      --name "lr${lr}_bs${bs}"
  done
done

# Run all queued experiments
dvc queue start --jobs 4

# Compare results
dvc exp show --only-changed

# Apply best experiment
dvc exp apply exp-best-id
git add params.yaml dvc.lock
git commit -m "Apply best hyperparameters"
```

# 11 Glossary

**Cache**
Local storage directory (`.dvc/cache`) where DVC stores actual file contents, indexed by hash.

**Checkpoint**
Special type of output that can be saved incrementally during training, useful for long-running experiments and allowing resumption.

**DAG**
Directed Acyclic Graph representing the pipeline's stage dependencies and execution order. Ensures no circular dependencies.

**Dependency**
File or directory that a stage depends on. Changes to dependencies trigger stage re-execution during `dvc repro`.

**DVC File**
Files with `.dvc` extension containing metadata and pointers to actual data stored in cache or remote storage.

**dvc.lock**
Lock file containing exact checksums and dependency versions, ensuring reproducibility across different environments.

**dvc.yaml**
Central pipeline definition file describing stages, dependencies, outputs, parameters, metrics, and plots.

**Experiment**
A variation of pipeline execution with different parameters, code, or data, tracked by DVC for comparison.

**Frozen Stage**
Stage marked with `frozen:    true` that won't execute automatically during `dvc repro`.

**Metric**
Numeric evaluation measure (accuracy, loss, F1) stored in JSON/YAML/CSV format for performance tracking.

**Output**
File or directory produced by a stage, tracked and cached by DVC for efficiency and reproducibility.

**Parameter**
Configuration value stored in `params.yaml`, tracked for experiment comparison and version control.

**Pipeline**
Series of connected stages forming a complete ML workflow, defined in `dvc.yaml`.

**Plot**
Data file (CSV/JSON) used for generating visualizations like learning curves, ROC curves, or confusion matrices.

**Remote Storage**
Cloud or network storage location (S3, GCS, Azure, SSH) for sharing DVC-tracked data with team members.

**Repro**
Short for "reproduce" - the `dvc repro` command that runs pipeline stages based on their dependencies.

**Stage**
Single step in DVC pipeline with defined command, dependencies, outputs, and optional parameters/metrics.

**Workspace**
Project's working directory containing checked-out files, source code, and pipeline configuration.

# 12   Conclusion

## 12.1   Key Takeaways

> **DVC Provides**:
>
> - **Version Control**: Git-like interface for data and models
>
> - **Reproducibility**: Guaranteed consistent results across runs
>
> - **Experimentation**: Track and compare parameter variations
>
> - **Collaboration**: Share data and pipelines with teams
>
> - **Scalability**: Handle datasets of any size
>
> - **Flexibility**: Works with any ML framework and cloud provider

## 12.2   Getting Started Checklist

1. **Initialize**: Set up Git and DVC in your project

2. **Configure Remote**: Add cloud storage for data sharing

3. **Track Data**: Use `dvc add` for large files

4. **Define Pipeline**: Create `dvc.yaml` with stages

5. **Set Parameters**: Configure `params.yaml`

6. **Run Pipeline**: Execute with `dvc repro`

7. **Track Metrics**: Monitor performance over time

8. **Version Control**: Commit `dvc.yaml` and `dvc.lock` to Git

9. **Share**: Push data with `dvc push`

10. **Collaborate**: Pull team changes with `git pull` and `dvc pull`

## 12.3   Best Practices Recap

> **Remember**:
> - Track ALL dependencies (data, code, configs)
> - Use `params.yaml` for all configurable values
> - Never cache metrics or plots
> - Set random seeds for reproducibility
> - Commit `dvc.yaml` and `dvc.lock` to Git
> - Push data regularly to remote storage
> - Document your pipeline for team understanding
> - Use descriptive stage and parameter names

## 12.4   Next Steps

1. **Practice**: Start with a small project

2. **Experiment**: Try different pipeline structures

3. **Integrate**: Add CI/CD automation

4. **Explore**: Use experiment tracking features

5. **Optimize**: Fine-tune performance and storage

6. **Share**: Collaborate with your team

## 12.5   Additional Resources

- **Documentation**: `https://dvc.org/doc`

- **Tutorial**: `https://dvc.org/doc/start`

- **GitHub**: `https://github.com/iterative/dvc`

- **Community**: `https://discuss.dvc.org`

- **Blog**: `https://dvc.org/blog`

- **YouTube**: `https://www.youtube.com/c/DVCorg`

- **Examples**: `https://github.com/iterative/example-get-started`

---

## Thank You!

---

*End of DVC Pipeline Configuration Guide*

Created by Sujil S — sujil9480@gmail.com