

# Pipeline Guide - LOE Forecasting System

---

## Quick Start

```
cd Main_project

# Run full pipeline (uses config.py settings)
python scripts/run_pipeline.py
```

---

## Configuration ( `src/config.py` )

All pipeline settings are controlled here. **No CLI arguments needed.**

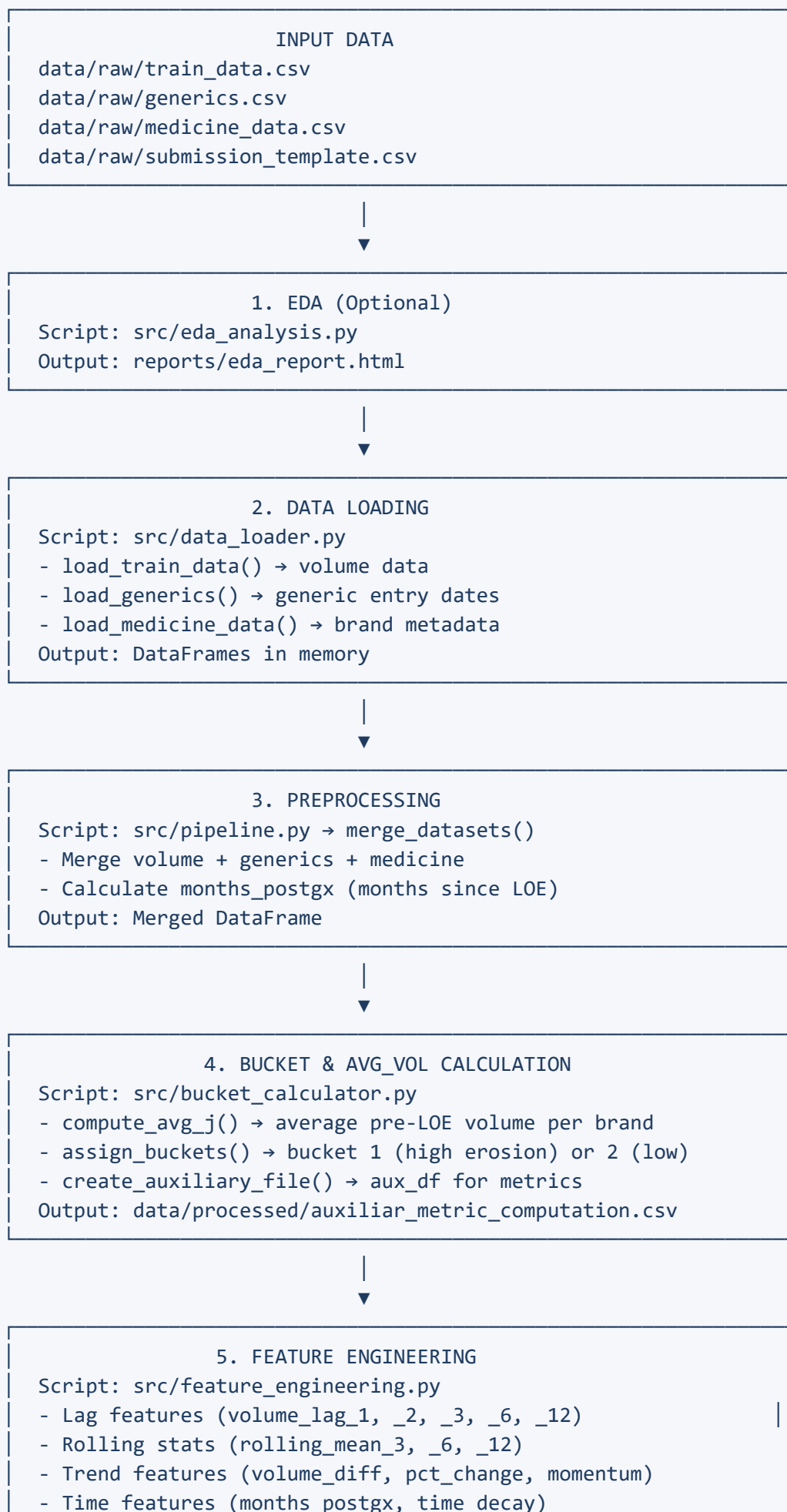
Setting	Options	Description
<code>TRAIN_MODE</code>	<code>"separate"</code> / <code>"unified"</code>	Train S1/S2 separately or single unified model
<code>TEST_MODE</code>	<code>True</code> / <code>False</code>	Use subset of brands for quick testing
<code>RUN_SCENARIO</code>	<code>"both"</code> / <code>"s1"</code> / <code>"s2"</code>	Which scenarios to run
<code>RUN_EDA</code>	<code>True</code> / <code>False</code>	Run exploratory data analysis
<code>RUN_TRAINING</code>	<code>True</code> / <code>False</code>	Train models
<code>RUN_SUBMISSION</code>	<code>True</code> / <code>False</code>	Generate submission file
<code>RUN_VALIDATION</code>	<code>True</code> / <code>False</code>	Validate submission format

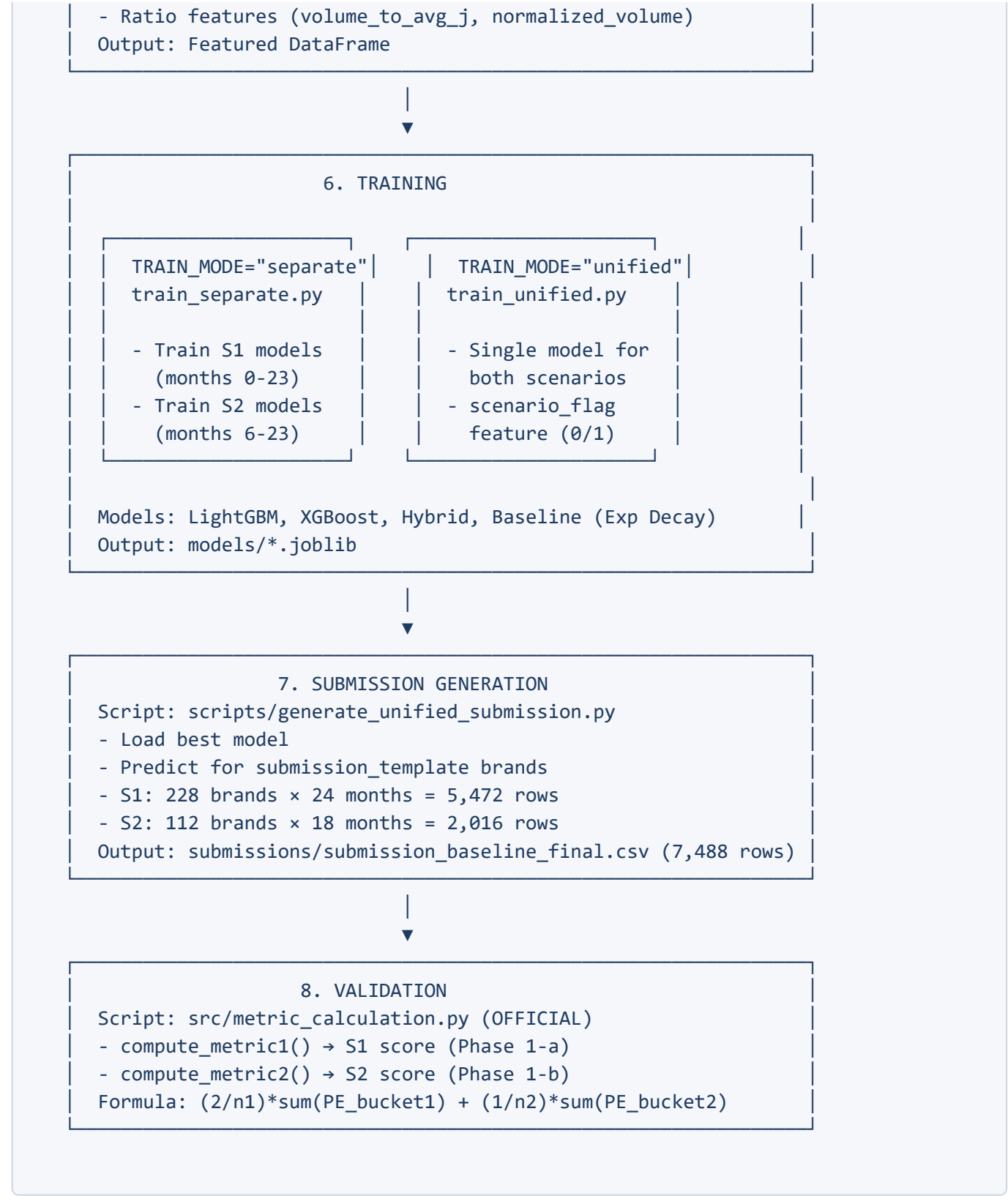
## Models Enabled

```
MODELS_ENABLED = {
    'baseline_exp_decay': True,
    'lightgbm': True,
    'xgboost': True,
    'hybrid_lightgbm': True,
    'arima': False, # Slow
}
```

---

## Pipeline Flow





File Reference

Scripts ( scripts/ )

File	Purpose	Run Command
run_pipeline.py	Main orchestrator - runs full pipeline	python scripts/run_pipeline.py

File	Purpose	Run Command
<code>train_models.py</code>	Training entry point (calls <code>train_separate</code> or <code>train_unified</code> )	Called by <code>run_pipeline.py</code>
<code>generate_unified_submission.py</code>	Generate submission CSV	Called by <code>run_pipeline.py</code>
<code>run_full_pipeline.py</code>	Alternative orchestrator	<code>python scripts/run_full_pipeline.py</code>

Source Modules ( `src/` )

File	Purpose
<code>config.py</code>	<b>Central configuration</b> - all settings here
<code>data_loader.py</code>	Load CSV files from <code>data/raw/</code>
<code>pipeline.py</code>	Merge datasets, coordinate preprocessing
<code>bucket_calculator.py</code>	Calculate <code>avg_vol</code> , assign buckets (1=high erosion, 2=low)
<code>feature_engineering.py</code>	Create ML features (lags, rolling, trends)
<code>models.py</code>	Model classes: <code>GradientBoostingModel</code> , <code>BaselineModels</code> , <code>HybridPhysicsMLModel</code>
<code>evaluation.py</code>	Cross-validation, model comparison
<code>metric_calculation.py</code>	<b>OFFICIAL</b> competition metrics
<code>eda_analysis.py</code>	Exploratory data analysis

Training Modules ( `src/training/` )

File	Purpose
<code>train_separate.py</code>	Separate training: S1 and S2 trained independently
<code>train_unified.py</code>	Unified training: single model with <code>scenario_flag</code> feature

Scenario Definitions ( `src/scenarios/` )

File	Purpose
<code>scenarios.py</code>	Centralized scenario definitions (S1: months 0-23, S2: months 6-23)

Inputs

File	Location	Description
------	----------	-------------

File	Location	Description
train_data.csv	data/raw/	Historical volume data (country, brand, date, volume)
generics.csv	data/raw/	Generic entry dates per brand
medicine_data.csv	data/raw/	Brand metadata
submission_template.csv	data/raw/	Template with brands/months to predict

Outputs

File	Location	Description
auxiliar_metric_computation.csv	data/processed/	avg_vol, bucket per brand
*.joblib	models/	Trained model files
model_comparison_*.csv	reports/	Training results comparison
submission_baseline_final.csv	submissions/	<b>Final submission file</b>
eda_report.html	reports/	EDA visualizations

Competition Metrics

Scenario 1 (Phase 1-a) - Zero Actuals

- Predict months 0-23 post-LOE
- No actual data given at prediction time
- PE formula includes terms for months 0-5, 6-11, 12-23

Scenario 2 (Phase 1-b) - Six Actuals

- Predict months 6-23 post-LOE
- Given actual data for months 0-5
- PE formula includes terms for months 6-11, 12-23

Final Score Formula

$$\text{Score} = (2/n1) \times \Sigma(\text{PE\_bucket1}) + (1/n2) \times \Sigma(\text{PE\_bucket2})$$

- Bucket 1: High erosion brands (mean\_ratio ≤ 0.25) - **weighted 2x**
- Bucket 2: Low erosion brands - **weighted 1x**
- Lower score = better

Example Workflow

```
# 1. Edit config.py as needed
#   Set TRAIN_MODE = "separate" or "unified"
#   Set TEST_MODE = False for full run

# 2. Run pipeline
cd Main_project
python scripts/run_pipeline.py

# 3. Check outputs
#   - models/ for trained models
#   - reports/ for comparison CSVs
#   - submissions/ for final submission
```

# Training Modes Explained

## Separate Mode ( TRAIN\_MODE = "separate" )

- Trains **two independent models**
- S1 model: trained on months 0-23, predicts with no actuals
- S2 model: trained on months 6-23, uses early months as features
- **Pros:** Specialized models for each scenario
- **Cons:** Less training data per model

## Unified Mode ( TRAIN\_MODE = "unified" )

- Trains **one model** for both scenarios
- Uses `scenario_flag` feature (0=S1, 1=S2)
- S2 samples include early post-LOE summary features
- **Pros:** More training data, shared learning
- **Cons:** May not specialize as well

# Troubleshooting

Issue	Solution
Import errors	Run from <code>Main_project/</code> directory
Missing data	Check <code>data/raw/</code> has all CSV files
Memory issues	Set <code>TEST_MODE = True</code> in config.py
Slow training	Disable <code>arima</code> in MODELS_ENABLED