

Novartis Datathon 2025 - Comprehensive TODO List

Competition Goal: Forecast generic erosion (normalized volume) for pharmaceutical drugs

Deadline: Phase 1A/1B submission deadline (check competition site)

Hero Model: CatBoost with scenario-specific sample weights

⚠️ When in doubt, complete Priority A items (CatBoost + features + validation + submission) before touching Priority B.

💻 Implementation Phases

Phase	Focus	Sections
Phase 1 (before first stable LB submission)	Core pipeline, baseline model, first submission	0–7 (Priority A only), 9.1, 9.4, 10.1–10.2, 11.1–11.3, 12.1–12.4
Phase 2 (after solid LB score)	Advanced features, unified logging, plotting, CV search, ensemble tuning	All remaining unchecked items (6.7–6.9, 8.* booster models, HPO, TabNet/FT-Transformer)

🔧 Canonical CLI Flags (Locked)

Flag	Format	Example
--scenario	Integer 1 or 2	--scenario 1
--model	Lowercase string	--model catboost
--split	train or test	--split train
--mode	train or test	--mode train
--data-config	Path to YAML	--data-config configs/data.yaml
--features-config	Path to YAML	--features-config configs/features.yaml
--run-config	Path to YAML	--run-config configs/run_defaults.yaml
--model-config	Path to YAML	--model-config configs/model_cat.yaml
--force-rebuild	Flag (no value)	--force-rebuild

Table of Contents

0. Design & Consistency Checks
1. Critical Path (Must Do)
2. Data Pipeline
3. Feature Engineering
4. Model Development
5. Training Pipeline

6. Validation & Evaluation
 7. Inference & Submission
 8. Experimentation & Optimization
 9. Testing & Quality Assurance
 10. Documentation & Presentation
 11. Colab/Production Readiness
 12. Competition Strategy
-

⌚ Recommended Implementation Order (Sections)

0 → 2 → 3 → 6 → 5 → 4 → 7 → 9 → 10 → 11 → 8 → 12

- **0 – Design & Consistency Checks:** configs, paths, CLI & docs alignment.
- **2 – Data Pipeline:** panels, caching, train/test behavior, leakage checks.
- **3 – Feature Engineering:** scenario-aware `make_features`, early-erosion features.
- **6 – Validation & Evaluation:** metrics, scenario simulation, official metric wrapper.
- **5 – Training Pipeline:** `train.py`, sample weights, run metadata, config wiring.
- **4 – Model Development:** baselines + CatBoost hero; other models after that.
- **7 – Inference & Submission:** `inference.py`, template alignment, aux file.
- **9 – Testing & QA:** unit + integration tests tightened around the full pipeline.
- **10 – Documentation & Presentation:** README, approach, methodology notes.
- **11 – Colab/Production Readiness:** Colab notebook, envs, performance tweaks.
- **8 – Experimentation & Optimization:** HPO, ensembles, ablations, smoothing.
- **12 – Competition Strategy:** LB management, final-week playbook, external-data checks.

Section 1 – Critical Path (Must Do) is a *running checklist* (EDA, first LB submission, score iterations) and should be followed in parallel as soon as the minimal 0→2→3→6→5→4 path is usable.

0. Design & Consistency Checks

Priority Note: Items directly tied to CatBoost + robust features + solid validation are **Priority A**; purely architectural items (full BaseModel hierarchy, advanced notebooks, stretch goals) are **Priority B / post-competition**.

0.0 Global Preparation & Repository Hygiene

- **Read and internalize documentation:**
 - Read `docs/planning/functionality.md` thoroughly
 - Read `docs/planning/approach.md`
 - Skim `docs/instructions/NOVARTIS_DATATHON_2025_COMPLETE_GUIDE_V2.md` and `docs/planning/question-set.md`
- **Ensure directory structure matches the spec:**
 - Verify top-level folders exist: `data/`, `src/`, `configs/`, `docs/`, `notebooks/`, `submissions/`, `tests/`, `artifacts/`
 - Verify `docs/guide/metric_calculation.py` exists
 - Verify `docs/guide/submitter_template.csv` exists
 - Verify `docs/guide/auxiliar_metric_computation_example.csv` exists
- **Verify environment files exist and are consistent:**

- `requirements.txt`, `env/colab_requirements.txt`, and/or `environment.yml` exist
- Ensure requirements list all needed packages and are tested (added pyarrow, pytest; all 92 tests pass)

0.1 Code-Documentation Alignment

- **Cross-check all items marked as ([x])** against the actual codebase and revert them to ([]) if implementation is partial or differs from `functionality.md` / `README.md`
- **Ensure function signatures and behaviors** in code match the descriptions in `docs/functionality.md` and `README.md` (e.g. `make_features`, `compute_pre_entry_stats`, `compute_metric1`, `generate_submission`, CLI entrypoints)
- **Align CLI examples** in `README.md`, `TODO.md` (Quick Commands) and actual `src/train.py` / `src/inference.py` argument names (`--data-config`, `--features-config`, `--run-config`, `--model-config`, `--scenario`, `--model` etc.)
 - Updated `README.md` CLI examples to use `--scenario 1` and `--scenario 2` (integer scenarios)
 - Updated `src/train.py` CLI to accept integer scenarios
 - Updated `src/evaluate.py` examples in `README` to match actual function signatures
- **Confirm all referenced paths exist** and are correctly used in code:
 - `docs/guide/metric_calculation.py`
 - `docs/guide/submission_template.csv`
 - `configs/*.yaml` (all referenced config files)

0.2 Configuration System (`configs/`)

- **Verify `configs/data.yaml` contains:**
 - `paths.raw_dir`, `paths.interim_dir`, `paths.processed_dir`, `paths.artifacts_dir`
 - `files.train`, `files.test`, `files.aux_metric`, `files.submission_template`
 - `columns` section: `id_keys`, `time_key`, `calendar_month`, `raw_target`, `model_target`, `meta_cols`
 - Added `mean_erotion` to `meta_cols` list
- **Verify `configs/features.yaml` contains:**
 - Sections: `pre_entry`, `time`, `generics`, `drug`, `scenario2_early`, `interactions`
 - Scenario-specific enable/disable flags for feature groups
- **Verify `configs/model_cat.yaml` contains:**
 - `model_type`, `params`, `training`, `categorical_features`
- **Verify `configs/run_defaults.yaml` contains:**
 - `reproducibility.seed`, `validation.*` settings
 - `scenarios.scenario1/scenario2` with `forecast_start`, `forecast_end`, `feature_cutoff`
 - `sample_weights` definition (time-window and bucket weights)
 - `logging.level`, `logging.log_to_file`
- **Configs as single source of truth:**

- Added `official_metric` section to `configs/run_defaults.yaml` with all metric weights and bucket threshold
- Updated `src/data.py compute_pre_entry_stats()` to accept `bucket_threshold` and `run_config` parameters
- Added `OFFICIAL_*` constants and `validate_config_matches_official()` to `src/evaluate.py`
- Documented that fallback metrics intentionally match official script (hardcoded values for exact correspondence)
- Added 5 new tests for config validation (92 tests total, all pass)

0.3 External Data & Constraints Rule-Check

- **Explicit rule-check on external data and constraints**
 - Re-read the latest official rules on:
 - Use of **external data**: Not explicitly prohibited but "Modeling Freedom: Any approach/model allowed"
 - Any restrictions on **model types** or **ensembles**: None - "explainability and simplicity valued"
 - Exact requirements for **reproducibility** and **code handover**: Required for finalists
 - Documented in `docs/planning/approach.md` (Appendix D.1):
 - **Decision: NOT using external data** - rationale documented (time investment, ROI uncertainty, complexity)
 - All constraints documented (no test target leakage, deterministic seeds, code reproducibility)

0.4 Code Fixes Applied (Implementation Log)

Summary: The following fixes were applied to achieve full Section 0 compliance.

- `src/data.py` :
 - Fixed syntax error: removed extra `"""` at line 1 (module docstring)
 - Fixed Unicode arrow character `←` in docstring (replaced with `<-`)
 - Added canonical constants: `ID_COLS`, `TIME_COL`, `CALENDAR_MONTH_COL`, `RAW_TARGET_COL`, `MODEL_TARGET_COL`
- `src/features.py` :
 - Changed `SCENARIO_CONFIG` keys from strings (`'scenario1'`, `'scenario2'`) to integers (`1`, `2`)
 - Added `_normalize_scenario()` helper to accept both int and string inputs for backward compatibility
 - Fixed month parsing: replaced `pd.to_datetime(df['month']).dt.month` with explicit month name mapping (`{'Jan': 1, ...}`)
 - Fixed FutureWarning in `groupby.apply()` by adding `include_groups=False`
- `src/train.py` :
 - Updated `--scenario` CLI argument to use `type=int`
 - Updated imports to use `META_COLS` from `src/data.py`
- `src/validation.py` :
 - Updated `simulate_scenario()` to use `_normalize_scenario()` from `features.py`
- `src/evaluate.py` :

- Updated `compute_bucket_metrics()` to handle integer scenarios
 - `src/inference.py` :
 - Updated `detect_test_scenarios()` to return `{1: [...], 2: [...]}` (integer keys)
 - `src/models/__init__.py` :
 - Implemented lazy imports for native-library models (LightGBM, XGBoost, CatBoost, NN) to prevent import failures when libraries are misconfigured
 - Linear models and baselines import eagerly (no native dependencies)
 - `configs/data.yaml` :
 - Added `mean_erosion` to `meta_cols` list
 - `tests/test_smoke.py` :
 - Complete rewrite to match actual function signatures
 - Updated to use integer scenarios (`scenario=1`, `scenario=2`)
 - Added `test_model_imports` using lazy imports (no skip needed)
 - Separated model interface test from import test
 - Added Section 0.2 config validation tests (5 new tests)
 - **Test results: 92 passed, 0 skipped, 0 warnings**
 - `tests/conftest.py` :
 - Created pytest configuration with warning filters
 - Filters torch/NumPy initialization warnings (external library)
 - Filters DataFrameGroupBy.apply FutureWarnings
 - `README.md` :
 - Updated CLI examples: `--scenario scenario1` → `--scenario 1`
 - Updated `src/evaluate.py` usage examples to match actual function signatures
-

1. Critical Path (Must Do)

1.1 Immediate Priorities (Day 1-2)

- **Verify data files are correctly placed** in `data/raw/TRAIN/` and `data/raw/TEST/`
- **Run smoke tests** to verify entire pipeline: `pytest tests/test_smoke.py -v` (372 passed, 0 skipped, 0 warnings)
- **Run EDA notebook** (`notebooks/00_eda.ipynb`) to understand data distributions
- **Establish baseline score** using `FlatBaseline` or `GlobalMeanBaseline`
- **Verify submission format** matches `docs/guide/submission_template.csv`

1.2 First Submission Target (Day 2-3)

- **Train CatBoost Scenario 1 model** with default hyperparameters (2025-11-29, RMSE=0.2488, Official=0.7692)
- **Train CatBoost Scenario 2 model** with default hyperparameters (2025-11-29, RMSE=0.2055, Official=0.2742)
- **Generate first submission file** for combined S1+S2 (submissions/submission_v1.csv)
- **Validate submission format** - 7488 rows, 340 series, months 0-23, no NaN/negative values
- **Submit to leaderboard** and record baseline score
- **Validate submission using** `metric_calculation.py` from docs/guide/

1.3 Score Improvement Iteration (Day 3-7)

- **Analyze feature importance** to identify key predictors
 - **Tune hyperparameters** using Optuna (focus on CatBoost)
 - **Implement ensemble** of top 2-3 models
 - **Add advanced features** based on EDA insights
 - **Re-submit** and track improvement
-

2. Data Pipeline

2.1 Data Loading (`src/data.py`)

- `load_raw_data()` - Load train/test CSV files
- `prepare_base_panel()` - Merge volume, generics, medicine info
- `compute_pre_entry_stats()` - Calculate avg_vol_12m, y_norm
- `handle_missing_values()` - Imputation strategies
- **Add validation** for expected column types and ranges
 - Added `validate_dataframe_schema()` for required columns
 - Added `validate_value_ranges()` with `VALUE_RANGES` dict
 - Added `validate_no_duplicates()` for key uniqueness
- **Add logging** for data loading statistics (rows, missing %)
 - Added `log_data_statistics()` helper function
 - Logs shape, missing value percentages, unique series count, time range
- **Handle edge cases:** Series with < 12 pre-entry months
 - Fallback 1: Use any available pre-entry months (`months_postgx < 0`)
 - Fallback 2: Use `ther_area` median
 - Fallback 3: Use global median (last resort)
 - Added `pre_entry_months_available` column for diagnostics
- **Add data caching** to speed up repeated loads (pickle/parquet)
 - Added `get_panel()` function with caching support
 - Parquet format if available, else pickle fallback
 - Added `clear_panel_cache()` utility
 - Added `_optimize_dtypes()` for storage efficiency
- **Define canonical constants** for id/time columns consistent with `configs/data.yaml`
 - `ID_COLS`, `TIME_COL`, `CALENDAR_MONTH_COL`, `RAW_TARGET_COL`, `MODEL_TARGET_COL`
 - `EXPECTED_DTYPES` and `VALUE_RANGES` for validation

2.2 Train vs Test Behavior for `compute_pre_entry_stats()`

- **Verify `is_train=True` behavior:** computes `y_norm`, `mean_erosion`, and `bucket`
- **Verify `is_train=False` behavior:** does **not** compute any target-dependent statistics (no `bucket`, no `mean_erosion`) and only computes `avg_vol_12m`
- **Implement robust fallback for `avg_vol_12m`** on test series with insufficient pre-entry history (e.g. global/`ther_area`-level averages)
 - 3-level fallback hierarchy: any pre-entry -> `ther_area` median -> global median
- **Document the train vs test behavior** in `docs/functionality.md`

- **Log distribution of avg_vol_12m** and bucket counts on train for sanity checking
 - Logs range, median, mean for avg_vol_12m
 - Logs bucket distribution with counts and percentages
- **Bucket definition:** Compute `bucket` as defined in `docs/functionality.md` and ensure code + configs implement the same rule
 - Bucket 1: mean_erosion <= 0.25 (high erosion)
 - Bucket 2: mean_erosion > 0.25 (low erosion)

2.3 META_COLS Consistency

- **Add a check that META_COLS** in `src/data.py` matches `columns.meta_cols` in `configs/data.yaml`
 - Added `verify_meta_cols_consistency()` function
 - Added test `test_meta_cols_consistency_check`
- **Update both files** if there is any mismatch (aligned in Section 0)

2.4 Data Leakage Audits

- **Implement systematic leakage audit** that confirms no feature uses:
 - Future `volume` values beyond the scenario cutoff
 - `bucket`, `mean_erosion`, or any other target-derived statistic
 - Test-set statistics (global means, encodings) computed using test data
- **Add automated leakage check** to run before training
 - `audit_data_leakage()` function in `src/data.py`
 - `run_pre_training_leakage_check()` function in `src/data.py`
 - `LEAKAGE_COLUMNS` constant for forbidden columns

2.5 Data Quality Checks

- **Verify no future leakage:** Features only use data before cutoff
 - `verify_no_future_leakage()` function in `src/data.py`
 - Called during `get_features()` to validate feature matrices
- **Check for duplicate rows** in panel
 - `validate_no_duplicates()` raises ValueError if duplicates detected
 - Called in `prepare_base_panel()`
- **Validate date continuity:** No gaps in months_postgx per series
 - `validate_date_continuity()` function in `src/data.py`
 - `get_series_month_coverage()` function for coverage statistics
- **Profile missing values** per column per scenario (via `log_data_statistics()`)
- **Check target distribution:** `y_norm` should be mostly 0-1.5
 - `VALUE_RANGES['y_norm'] = (0, 5)` with warnings for outliers

2.6 Data Splits

- **Implement time-based CV split** for more realistic validation
 - `create_temporal_cv_split()` function in `src/validation.py`
- **Stratify by bucket AND therapeutic area** for balanced folds
 - `create_validation_split()` supports `stratify_by` list of columns

- **Create holdout set** from training data for final validation
 - `create_holdout_set()` function in `src/validation.py`
- **Document split rationale** in approach.md

2.7 Utility Functions (`src/utils.py`)

- `set_seed(seed)` : Sets Python, NumPy (and torch if available) RNGs for reproducibility
- `setup_logging(level, log_file)` : Configures the root logger, avoids duplicate handlers
- `timer(name)` : Context manager that logs start/end and elapsed seconds
- `load_config(path)` : Loads YAML, raises clear errors if missing/invalid
- `get_path(config, key)` : Helper returning `Path` objects from config paths

2.8 Persisted Data Build (raw → interim → processed)

Goal: Have a **single, well-defined data build pipeline** that:

- Reads CSVs from `raw/`
- Builds cleaned panels and saves them to `interim/`
- Builds scenario-specific feature matrices and saves them to `processed/`
- **Without creating extra modules** (only extend `src/data.py`, `src/features.py`, and configs)

2.8.1 Align Paths in `configs/data.yaml`

- Set `paths.raw_dir` to the existing `./raw` directory (containing `TRAIN/` and `TEST/`)
- Set `paths.interim_dir` to `./interim`
- Set `paths.processed_dir` to `./processed`
- Ensure no new top-level `data/` folder is created implicitly (reuse the existing `raw/`, `interim/`, `processed/`)

2.8.2 Implement Cached Panel Builder in `src/data.py`

- Add function `get_panel(split: str, config, use_cache: bool = True, force_rebuild: bool = False) -> pd.DataFrame` that:
 - Determines the panel path, e.g. `interim/panel_{split}.parquet` (e.g. `panel_train.parquet`, `panel_test.parquet`)
 - If `use_cache=True` and the parquet file exists and `force_rebuild=False`, loads it and returns
 - Otherwise:
 - Calls `load_raw_data(config, split=split)` to read CSVs from `raw/TRAIN` or `raw/TEST`
 - Calls `prepare_base_panel(...)` and `handle_missing_values(...)`
 - Calls `compute_pre_entry_stats(..., is_train=(split == "train"))`
 - Saves the resulting panel to `interim/panel_{split}.parquet`
 - Returns the panel DataFrame
- Ensure this function is the only entry point used by training/inference to obtain panels, to avoid duplicated logic

2.8.3 Implement Persisted Feature Matrices in `src/features.py`

- Add function `get_features(split: str, scenario: int, mode: str, config, use_cache: bool = True, force_rebuild: bool = False) -> Tuple[pd.DataFrame, pd.Series, pd.DataFrame]` that:
 - Computes the output path based on `paths.processed_dir`, e.g.:
 - `processed/features_{split}_scenario{scenario}_{mode}.parquet`
 - `processed/target_{split}_scenario{scenario}_{mode}.parquet` (optional, for `y`)
 - `processed/meta_{split}_scenario{scenario}_{mode}.parquet` (optional, for `META_COLS`)
 - If `use_cache=True` and files exist and `force_rebuild=False`, loads and returns `(X, y, meta)` (for `mode="train"`; for `mode="test"` `y` is `None`)
 - Otherwise:
 - Uses `get_panel(split, config)` to obtain the base panel
 - Calls `make_features(panel_df, scenario=scenario, mode=mode, config=features_config)`
 - Splits into `(X, y, meta)` via `split_features_target_meta`
 - Saves `X` (features) to `processed/features_...parquet`
 - Saves `y` (if not None) and `meta` (if needed) similarly
 - Returns `(X, y, meta)`
- Add `clear_feature_cache()` function to clear cached feature files

2.8.4 Wire Training and Inference to Use Persisted Features

- In `src/train.py`, modify `train_scenario_model` / `run_experiment` to:
 - Call `get_features(split="train", scenario, mode="train", ...)` instead of manually building panels + features
 - Optionally support a flag `--force-rebuild` in the CLI that sets `force_rebuild=True` for both `get_panel` and `get_features`
 - Added `--no-cache` flag to disable feature caching
- In `src/inference.py`, modify the inference pipeline to:
 - Use `get_panel()` for test data loading
 - Added `--force-rebuild` CLI flag

2.8.5 Add Minimal CLI Entrypoints for Data Build (No New Modules)

- In `src/data.py`, add a small `if __name__ == "__main__":` block using `argparse` to:
 - Accept arguments like `--split train|test`, `--scenario 1|2`, `--mode train|test`, `--force-rebuild`
 - Accept `--data-config` and `--features-config` paths
 - Call `get_panel(...)` and/or `get_features(...)` to pre-build and cache data
- Add commands to `README.md` / `TODO.md` "Quick Commands" (see below)
- Emphasize that no new Python files (beyond existing modules) should be created for this; all logic stays inside `src/data.py` and `src/features.py`

2.8.6 Implement Basic Schema and Shape Validation at Each Stage

- After building each panel, assert that:
 - Required columns (`country`, `brand_name`, `months_postgx`, `volume`, `n_gxs`, drug attributes) are present
 - There are no duplicate keys (`country`, `brand_name`, `months_postgx`)
 - `validate_panel_schema()` function added to `src/data.py`
- After building each feature matrix:
 - Assert that feature columns contain **no META_COLS**
 - Assert that `X.shape[0] == len(y)` (for `mode="train"`)
 - Log shapes, e.g. `X_train_s1: (n_rows, n_features)`
 - `validate_feature_matrix()` function added to `src/features.py`

2.8.7 Optimize Storage Formats and Dtypes

- Use Parquet format for both panels and features to reduce disk usage and speed up I/O
 - In `prepare_base_panel` and `make_features`:
 - Cast high-cardinality categoricals (`country`, `brand_name`, `ther_area`, `main_package`) to `category` dtype before saving
 - Optionally downcast numeric types (`int64` → `int32`, `float64` → `float32`) when safe, and log the downcasting
 - `_optimize_dtypes()` function updated in `src/data.py`
 - Preserves precision for critical columns (`y_norm`, `volume`, `avg_vol_12m`)
-

3. Feature Engineering

3.0 Feature Mode Handling (`make_features`)

- Ensure `make_features(panel_df, scenario, mode)` behaves correctly:
 - Creates `y_norm` **only** when `mode="train"`
 - Never touches target-related columns for `mode="test"` (no `y_norm`, no `bucket`)
 - Respects scenario-specific cutoffs internally (`months_postgx < 0` for S1, `< 6` for S2)
- Ensure `configs/features.yaml` can enable/disable features per scenario (e.g. Scenario 2 early-erosion features are not accidentally computed for Scenario 1)
- Confirm `make_features` loads and applies scenario-specific settings from `features.yaml`, not hardcoded logic
 - Added `_load_feature_config()` helper to merge user config with defaults
- Implement helper functions driven by `configs/features.yaml`:
 - `add_pre_entry_features`, `add_time_features`, `add_generics_features`, `add_drug_features`, `add_early_erosion_features`
 - All functions now accept config dict parameter

3.1 Pre-Entry Features (`src/features.py`)

- `avg_vol_3m`, `avg_vol_6m`, `avg_vol_12m` - Rolling averages
- `log_avg_vol_12m` - Log-transformed scale normalization
- `pre_entry_trend` - Linear slope of pre-entry volume
- `pre_entry_volatility` - Coefficient of variation
- **Add pre_entry_max** - Maximum volume in pre-entry period

- **Add pre_entry_min** - Minimum volume in pre-entry period
- **Add volume_growth_rate** - $(\text{vol_3m} - \text{vol_12m}) / \text{vol_12m}$
- **Add pre_entry_trend_norm** - Normalized trend for scale invariance
- **Add pre_entry_max_ratio, pre_entry_min_ratio, pre_entry_range_ratio**
- **Add seasonal pattern detection** from pre-entry months
 - Implemented `_add_seasonal_features()` function with:
 - `seasonal_amplitude` : Max deviation from mean by month-of-year
 - `seasonal_peak_month` : Month with highest volume
 - `seasonal_trough_month` : Month with lowest volume
 - `seasonal_ratio` : Peak-to-trough ratio
 - `seasonal_q1_effect`, `seasonal_q2_effect`, `seasonal_q3_effect`, `seasonal_q4_effect` : Quarter-wise deviations
 - Documented in `configs/features.yaml`

3.2 Time Features

- `months_postgx` - Direct time index
- `months_postgx_squared` - Quadratic decay
- `months_postgx_cube` - Cubic for flexible decay
- `is_post_entry` - Binary flag
- `time_bucket` - Categorical (pre, early, mid, late)
- `is_early`, `is_mid`, `is_late` - One-hot time bucket encoding
- **Add time decay curve features:** $\exp(-\alpha * \text{months_postgx})$
 - `time_decay` ($\alpha=0.1$) and `time_decay_fast` ($\alpha=0.2$)
- **Add month_of_year** with cyclical encoding (`month_sin`, `month_cos`)
- **Add quarters** (Q1-Q4) for seasonality
 - `quarter` and `is_q1`, `is_q2`, `is_q3`, `is_q4` one-hot encoding
- **Add is_year_end** flag for December submissions
- **Add is_year_start** flag for January
- **Add sqrt_months_postgx** for slower decay

3.3 Generics Competition Features

- `n_gxs` - Current number of generics
- `has_generic` - Binary ($n_{gxs} > 0$)
- `multiple_generics` - Binary ($n_{gxs} \geq 2$)
- `many_generics` - Binary ($n_{gxs} \geq 5$)
- `n_gxs_pre_cutoff_max` - Maximum generics up to cutoff month
- **Add first_generic_month** - Month of first generic entry
- **Add months_since_first_generic** - Time since first competitor
- **Add had_generic_pre_entry** - Binary flag
- **Add generic_entry_speed** - Rate of new generic entries
- **Add log_n_gxs, log_n_gxs_at_entry** - Log transforms
- **Add n_gxs_bin** - Categorical binning (none/one/few/several/many)
- **Add expected_future_generics** - Future `n_gxs` is EXOGENOUS (provided in test data)
 - Implemented `_add_future_generics_features()` function with:
 - `n_gxs_at_month_12` : Expected `n_gxs` at month 12

- `n_gxs_at_month_23` : Expected n_gxs at end of forecast horizon
- `n_gxs_change_to_12`, `n_gxs_change_to_23` : Change features
- `n_gxs_max_forecast` : Maximum n_gxs over forecast horizon
- `expected_new_generics` : Number of new generics expected
- NOTE: This is NOT leakage - n_gxs is exogenous (from competitive intelligence)

3.4 Drug Characteristics (Static Features)

- `ther_area` - Therapeutic area (categorical)
- `main_package` - Dosage form (categorical)
- `hospital_rate` - Hospital percentage (0-100)
- `biological` - Boolean
- `small_molecule` - Boolean
- `hospital_rate_bin` - Binned hospital rate (low/medium/high)
- `hospital_rate_norm` - Normalized 0-1
- `is_hospital_drug`, `is_retail_drug` - Binary flags
- **Add `is_injection`** derived from `main_package`
- **Add `is_oral`** derived from `main_package`
- **Add `ther_area_encoded`, `main_package_encoded`, `country_encoded`** - Label encoding
- **Add `biological_x_n_gxs`** - Interaction feature
- **Add `hospital_rate_x_time`** - Interaction feature
- **Add `ther_area erosion prior`** - Historical avg erosion by area
 - Implemented `add_target_encoding_features()` with K-fold cross-validation
 - CRITICAL: Uses K-fold to compute encoding only from OTHER series (no leakage)
 - Smoothing parameter for regularization
- **Add `country_effect`** - Country-level erosion patterns
 - `country_erosion_prior` via same target encoding framework
- **Add `ther_area encoding`** (target encoding with K-fold)
 - Configurable via `configs/features.yaml` `target_encoding` section

3.5 Scenario 2 Specific Features

- `avg_vol_0_5` - Mean volume over months [0, 5]
- `erosion_0_5` - Early erosion signal
- `trend_0_5` - Linear slope over months 0-5
- `trend_0_5_norm` - Normalized trend
- `drop_month_0` - Initial volume drop
- **Add `month_0_to_3_change`** - Short-term erosion rate
- **Add `month_3_to_5_change`** - Medium-term erosion rate
- **Add `avg_vol_0_2, avg_vol_3_5`** - Window averages
- **Add `erosion_0_2, erosion_3_5`** - Window erosion ratios
- **Add `recovery_signal`** - If volume increases after initial drop
 - Defined: boolean flag if average volume in months [3-5] is higher than in months [0-2]
- **Add `recovery_magnitude`** - Size of recovery (clipped to [-1, 1])
- **Add `competition_response`** - n_gxs change in 0-5
 - Defined: change in `n_gxs` between months 0 and 5
- **Add `erosion_per_generic`** - Erosion per new generic competitor

3.6 Feature Scenario & Cutoff Validation

- **Verify that any feature using post-entry volume** (e.g. `avg_vol_0_5`, `erosion_0_5`, `trend_0_5`) is only used:
 - In Scenario 2
 - And only derived from months strictly before the scenario's forecast start (6–23 for S2)
- **Add automated tests** to enforce these constraints
 - `validate_feature_leakage()` function
 - `validate_feature_cutoffs()` function
 - Test `test_feature_leakage_validation`
 - Test `test_scenario1_no_early_erosion_features`

3.7 Interaction Features

- **Implement interaction feature framework**
 - `add_interaction_features()` function
 - Configurable via `features.yaml` interactions section
- **Add `n_gxs` × `biological`** - Competition impact on biologics (via `biological_x_n_gxs`)
- **Add `hospital_rate` × `months_postgx`** - Time-hospital interaction (via `hospital_rate_x_time`)
- **Add `ther_area` × `erosion_pattern`** - Area-specific curves
 - `ther_area_x_early_erosion` : Interaction with early erosion (S2 only)
 - `ther_area_erosion_x_time` : Interaction with time
 - Created automatically when `target_encoding` is enabled

3.8 Feature Selection

- **Analyze feature correlations** and remove redundant features
 - `analyze_feature_correlations()` function
 - `remove_redundant_features()` function
- **Use permutation importance** post-training
 - `compute_feature_importance_permutation()` function
 - `select_features_by_importance()` helper
- **Implement feature summary** for analysis
 - `get_feature_summary()` function
- **Try recursive feature elimination** for linear models (optional, use if needed)
- **Document final feature set** for each scenario (to be done during experimentation)

3.9 Code Additions (Implementation Log)

Summary: The following code additions implement Section 3.

- `src/features.py` :
 - Added `FORBIDDEN_FEATURES` constant for leakage prevention
 - Added `_load_feature_config()` to merge config with defaults
 - Updated `make_features()` to use config-driven feature generation
 - Updated `add_pre_entry_features()` with config support and new features
 - Added `_add_seasonal_features()` for seasonal pattern detection (Section 3.1)

- Updated `add_time_features()` with decay curves, quarters, year-end flags
- Updated `add_generics_features()` with entry speed, binning, log transforms
- Added `_add_future_generics_features()` for exogenous future n_gxs (Section 3.3)
- Updated `add_drug_features()` with is_injection, is_oral, interactions
- Added `add_target_encoding_features()` with K-fold cross-validation (Section 3.4)
- Updated `add_early_erosion_features()` with change windows, recovery, competition
- Added `add_interaction_features()` for configurable interactions
- Added `validate_feature_leakage()` for leakage detection
- Added `validate_feature_cutoffs()` for cutoff validation
- Added `split_features_target_meta()` for clean feature/target separation
- Added `get_categorical_feature_names()` and `get_numeric_feature_names()` helpers
- Added Feature Selection utilities (Section 3.8):
 - `analyze_feature_correlations()`
 - `compute_feature_importance_permutation()`
 - `select_features_by_importance()`
 - `get_feature_summary()`
 - `remove_redundant_features()`
- `tests/test_smoke.py`:
 - Added `test_feature_config_loading`
 - Added `test_pre_entry_features`
 - Added `test_time_features`
 - Added `test_generics_features`
 - Added `test_scenario2_early_erosion_features`
 - Added `test_feature_leakage_validation`
 - Added `test_scenario1_no_early_erosion_features`
 - Added `test_make_features_mode_train_vs_test`
 - Added `test_split_features_target_meta`
 - Added `test_interaction_features`
 - Added Section 3 tests:
 - `TestSeasonalFeatures` : `test_seasonal_features_created`, `test_seasonal_amplitude_captures_pattern`
 - `TestFutureGenericsFeatures` : `test_future_generics_features_created`, `test_future_generics_values_correct`
 - `TestTargetEncodingFeatures` : `test_target_encoding_function_exists`, `test_target_encoding_creates_features`
 - `TestFeatureSelection` : `test_correlation_analysis`, `test_feature_importance`, `test_feature_summary_function`, `test_remove_redundant_features`
 - `TestInteractionFeatures` : `test_ther_area_erosion_interaction_created`
- **Test results: 104 passed, 0 skipped, 0 warnings**

4. Model Development

Section 4 Status: All Priority A infrastructure is COMPLETE

- BaseModel interface implemented and verified
- All baseline models implemented (Flat, GlobalMean, Trend, HistoricalCurve)
- All linear models implemented with polynomial features
- CatBoost, LightGBM, XGBoost, NN models implemented
- All 4 ensemble methods implemented (Averaging, Weighted, Stacking, Blending)
- Model factory functions working in train.py
- **198 tests passing, 0 skipped, 0 warnings**

Remaining items in Section 4 are:

- **Hyperparameter tuning** → Section 8 (Experimentation & Optimization)
- **Recording baseline scores** → Requires running experiments
- **Segment-specific models** (4.10, 4.11) → Experimentation tasks for Section 8
- **TabNet/FT-Transformer** → Priority B / Nice-to-Have

4.0 Priority Classification & Model Interface

Priority A (Must-Have): CatBoost + robust features + solid validation + simple ensemble

Priority B (Nice-to-Have): Advanced neural architectures, complex augmentation, extensive HPO

- **Ensure all models implement a common interface** (`fit`, `predict`, `save`, `load`, optional `get_feature_importance`) so that `train_scenario_model` and `inference` can treat them polymorphically
 - All model classes (CatBoost, LightGBM, XGBoost, Linear, NN, baselines, ensembles) implement BaseModel interface
 - `get_model_class()` factory function in `src/models/__init__.py`
 - `_get_model()` function in `src/train.py` for dynamic model instantiation
- **Verify BaseModel abstract class** in `src/models/base.py` defines the required interface
 - Abstract methods: `fit`, `predict`, `save`, `load`
 - Optional method: `get_feature_importance`

4.1 Baseline Models (`src/models/linear.py`)

- `FlatBaseline` - Always predicts 1.0 (no erosion)
- `GlobalMeanBaseline` - Predicts mean erosion curve
- `TrendBaseline` - Extrapolates pre-entry trend
- **Implement HistoricalCurveBaseline** - Matches to similar historical series (2025-01-XX)
 - Uses K-nearest neighbors to find similar series based on pre-entry features
 - Matching features: ther_area, hospital_rate, biological, small_molecule, log_avg_vol_12m, pre_entry_trend, pre_entry_volatility, n_gxs
 - Uses averaged erosion curves from similar series as predictions
 - Configurable: n_neighbors, metric, weights
- **Record baseline scores** for both scenarios

4.2 Linear Models (`src/models/linear.py`)

- `Ridge` - L2 regularized linear regression
- `Lasso` - L1 regularized (sparse)
- `ElasticNet` - Combined L1+L2

- `HuberRegressor` - Robust to outliers
- **Tune regularization strength** (alpha) via CV
- **Try polynomial features** (degree 2) with linear models (2025-01-XX)
 - `LinearModel` supports `use_polynomial: true` and `polynomial_degree` config options
 - Uses `sklearn PolynomialFeatures` for feature expansion
- **Use linear model coefficients** for interpretability insights

4.3 CatBoost (`src/models/cat_model.py`) - Hero Model

- Basic implementation with native categorical support
- Sample weight support via Pool
- Early stopping support
- Feature importance extraction
- **Tune depth** (try 4, 6, 8)
- **Tune learning_rate** (try 0.01, 0.03, 0.05)
- **Tune l2_leaf_reg** (try 1, 3, 5, 10)
- **Try loss_function = 'MAE'** instead of RMSE
- **Try loss_function = 'Quantile'** for different quantiles
- **Implement custom objective** aligned with official metric

4.4 LightGBM (`src/models/lgbm_model.py`)

- Basic implementation
- Sample weight support
- Early stopping
- **Tune num_leaves** (try 15, 31, 63)
- **Tune min_data_in_leaf** (try 10, 20, 50)
- **Try dart boosting** for regularization
- **Compare speed vs CatBoost**

4.5 XGBoost (`src/models/xgb_model.py`)

- Basic implementation with DMatrix
- Sample weight support
- Early stopping
- **Tune max_depth** (try 4, 6, 8)
- **Tune min_child_weight** (try 1, 5, 10)
- **Try colsample_bytree** (try 0.6, 0.8, 1.0)
- **Enable GPU** if available (`tree_method='gpu_hist'`)

4.6 Neural Network (`src/models/nn.py`)

- SimpleMLP with configurable layers
- Dropout and batch normalization
- Early stopping
- WeightedRandomSampler for sample weights
- **Experiment with architecture** (try [128, 64], [512, 256, 128])

- **Try different activations** (GELU, SiLU instead of ReLU)
- **Add residual connections** for deeper networks
- **Try embedding layers** for categorical features

4.7 Neural Network - Nice-to-Have / Post-Competition

Note: These are lower priority and should not distract from core datathon-critical work

- **Implement TabNet** architecture
- **Implement FT-Transformer** architecture

4.8 Ensemble Methods (Priority A)

- **Implement simple averaging** of predictions (2025-01-XX)
 - `AveragingEnsemble` class in `src/models/ensemble.py`
 - Simple mean of base model predictions
 - Configurable prediction clipping
- **Implement weighted averaging** (tune weights on validation) (2025-01-XX)
 - `WeightedAveragingEnsemble` class in `src/models/ensemble.py`
 - Weight optimization via `scipy.optimize.minimize`
 - Supports MSE or MAE optimization metric
 - Weights are constrained to be non-negative and sum to 1
- **Implement stacking** with meta-learner (2025-01-XX)
 - `StackingEnsemble` class in `src/models/ensemble.py`
 - K-fold cross-validation for OOF predictions
 - Default meta-learner: Ridge regression
 - Option to include original features in meta-learner
- **Implement blending** with hold-out predictions (2025-01-XX)
 - `BlendingEnsemble` class in `src/models/ensemble.py`
 - Holdout fraction configurable (default: 0.2)
 - Meta-learner trained on holdout predictions
- **Factory function** `create_ensemble(models, method, **kwargs)` for easy ensemble creation
- **Try CatBoost + LightGBM + XGBoost ensemble**
- **Document ensemble weights** for reproducibility

4.9 Code Additions (Implementation Log)

Summary: The following code additions implement Section 4.

- `src/models/__init__.py` :
 - Added exports for `HistoricalCurveBaseline`
 - Added exports for ensemble models: `AveragingEnsemble`, `WeightedAveragingEnsemble`, `StackingEnsemble`, `BlendingEnsemble`, `create_ensemble`
 - Updated `get_model_class()` mapping with all new model types
- `src/models/linear.py` :
 - Added `HistoricalCurveBaseline` class (~200 lines)
 - Added polynomial feature support to `LinearModel`

- `src/models/ensemble.py` (NEW FILE):
 - `AveragingEnsemble` - Simple mean ensemble
 - `WeightedAveragingEnsemble` - Optimized weight ensemble
 - `StackingEnsemble` - Two-level stacking with meta-learner
 - `BlendingEnsemble` - Holdout blending ensemble
 - `create_ensemble()` factory function
 - All ensembles implement `BaseModel` interface (fit, predict, save, load, `get_feature_importance`)
- `src/train.py` :
 - Updated `_get_model()` function with all new model types
 - Added: historical_curve, knn_curve, averaging, weighted, stacking, blending, trend, nn
- `tests/test_smoke.py` :
 - `TestBaseModelInterface` - Tests for `BaseModel` ABC and factory function
 - `TestHistoricalCurveBaseline` - Tests for historical curve baseline
 - `TestLinearModelPolynomial` - Tests for polynomial features
 - `TestEnsembleModels` - Tests for all 4 ensemble types
 - `TestTrainGetModel` - Tests for `_get_model` function
 - `TestTrendBaseline` - Tests for trend baseline save/load
 - **Test results: 198 passed, 0 skipped, 0 warnings**

4.10 Segment-Specific Models for High-Impact Regions (Priority A - Experimentation)

Rationale: The metric heavily weights bucket 1 and early months. Dedicated models for these segments can be decisive. **Note:** These are experimentation tasks that require running actual training experiments with the full pipeline. Implementation of infrastructure (model classes, ensemble methods) is complete. The experiments should be run as part of Section 8 (Experimentation & Optimization).

- **Scenario- and segment-specific models for high-impact regions**
 - Train **separate CatBoost models for bucket 1 vs bucket 2** (using bucket only on train, never as feature). At inference, assign each test series to a "pseudo-bucket" using pre-entry features + early-post-entry dynamics (Scenario 2) via a small classifier, and route its prediction to the corresponding expert model.
 - Train **early-window-focused models**:
 - For Scenario 1: a dedicated model optimised only on months 0–5 targets (with higher sample weights), then blend its predictions with the full-horizon model.
 - For Scenario 2: a dedicated model focused on months 6–11 (the heavy-weight window), blended with the full-horizon one.
 - Evaluate whether **segment-specific models + blending** beat a single global CatBoost in CV and on the leaderboard proxy.

4.11 Hierarchical / Segmented Modelling by Country & Therapeutic Area (Priority B - Experimentation)

Note: These are advanced experimentation tasks. The infrastructure (ensemble methods, model factory) is in place. Experiments should be run as part of Section 8.

- **Hierarchical / segmented modelling by country and therapeutic area**

- Cluster series into **homogeneous groups** (e.g. by country, therapeutic area, hospital_rate_bin, biologic vs small molecule).
 - For each group with sufficient data, train:
 - A **group-specific CatBoost model**.
 - Or a **shared global model + group-specific bias adjustment** (e.g., post-hoc calibration per group).
 - Compare:
 - Global-only model vs segmented ensemble on the unified CV scheme.
 - Pay special attention to segments that contribute most to the metric (e.g., high-erosion bucket 1 series in large markets).
 - If segmented models improve the metric, integrate them into the main **ensemble/blending step** and document the strategy.
-

5. Training Pipeline

5.1 Core Training (`src/train.py`)

- `split_features_target_meta()` - Separate columns
- `compute_sample_weights()` - Time + bucket weights
- `train_scenario_model()` - Train single model
- `run_experiment()` - Full experiment loop
- CLI interface with argparse
- **Add cross-validation loop** (`run_cross_validation()` K-fold training with series-level splits)
- **Add OOF prediction saving** (OOF predictions saved in `run_cross_validation()`)
- **Add experiment tracking** (MLflow/W&B integration) (2025-01-XX)
 - `init_experiment_tracking()` - Initialize MLflow or W&B with graceful fallback
 - `log_metrics_to_tracker()` - Log metrics to active tracker
 - `ExperimentTracker` class for unified tracking API
 - `MLFLOW_AVAILABLE`, `WANDB_AVAILABLE` constants for optional dependencies
 - `configs/run_defaults.yaml` updated with `experiment.tracking` config section
- **Add checkpoint saving** for resume training (2025-01-XX)
 - `save_model_checkpoint()` - Save model with metadata (epoch, metrics, config)
 - `load_model_checkpoint()` - Resume from checkpoint
 - `TrainingCheckpoint` class for checkpoint management
 - CLI options `--save-checkpoint`, `--load-checkpoint`
- **Add config hashing** for reproducibility (2025-01-XX)
 - `compute_config_hash()` - Deterministic hash of config dictionaries
 - Save config hash with experiment artifacts in `run_experiment()`
 - Config hash saved to `artifacts/run_id/config_hash.txt`
 - `save_config_snapshot()` - Save exact copies of all config files to `artifacts/run_id/configs/`

5.2 CLI Consistency & Help

- **Ensure `src/train.py` and `src/inference.py` implement `--help`** with clear argument documentation

- **Update TODO.md Quick Commands** to use exactly the same argument names as the code
(already consistent: `--data-config`, `--features-config`, `--run-config`, `--model-config`)

5.3 Experiment Metadata Logging

- **Log at the start of each training run:**
 - Scenario, model type, config paths
 - Random seed
 - Git commit hash (if available) - `get_git_commit_hash()`
 - Dataset sizes (number of series, number of rows) - `get_experiment_metadata()`
 - Config hash for reproducibility (`config_hash` in `metadata.json`)
- **Save a small JSON/YAML metadata file** in `artifacts/` per run (`metadata.json`, `config_snapshot.yaml`, `metrics.json` in `run_experiment()`)

5.4 Sample Weights Refinement

- Time-window weights (50/20/10 for S1, 50/30/20 for S2)
- Bucket weights (2× for bucket 1)
- **Fine-tune time weights** based on official metric formula (2025-01-XX)
- **Add month-level weights** for 20% monthly component (2025-01-XX)
- **Experiment with sqrt/log transformations** of weights (2025-01-XX)
 - `compute_sample_weights(weight_transform=...)` - Apply sqrt/log/rank transforms
 - `transform_weights()` function in `src/train.py`
 - CLI option `--weight-transform` (identity/sqrt/log/rank)
- **Validate weights** correlate with metric improvement (2025-01-XX)
 - `validate_weights_metric_alignment()` - Compute correlation between weights and metric improvement
 - `validate_weights_correlation()` helper function
- **Set weights according to** `configs/run_defaults.yaml` (`compute_sample_weights()` reads from config)
- **Explicit alignment of training loss with official metric** (2025-01-XX)
 - Derive **exact per-row weights** that reproduce the official metric contribution:
 - Early windows vs rest of horizon
 - Bucket 1 vs bucket 2
 - Monthly 20% component
 - Implement a "**metric-aligned weight calculator**"
(`compute_metric_aligned_weights()`)
 - Scenario 1: $0.2 \times \text{monthly} + 0.5 \times (0-5) + 0.2 \times (6-11) + 0.1 \times (12-23)$
 - Scenario 2: $0.2 \times \text{monthly} + 0.5 \times (6-11) + 0.3 \times (12-23)$
 - Bucket weighting: 2× for bucket 1
 - Inverse avg_vol weighting for smaller series
 - Integrate via `compute_sample_weights(use_metric_aligned=True)`
 - Plug these weights into CatBoost/GBMs and validate on a small toy example that the **weighted RMSE at row level aggregates to the same series-level metric** (up to numerical noise).

5.5 Hyperparameter Optimization

- **Implement Optuna integration** for CatBoost (2025-01-XX)
 - `run_hyperparameter_optimization()` - Main Optuna optimization function
 - `create_optuna_objective()` - Create objective function for Optuna
 - `OPTUNA_AVAILABLE` constant for graceful fallback
 - CatBoostPruningCallback integration for early termination
- **Define search space** in configs/model_cat.yaml (2025-01-XX)
 - Default HPO search space with parameter ranges
- **Use pruning** for early termination of bad trials
- CLI options `--hpo-trials`, `--hpo-timeout` for optimization control
- **Run for 100+ trials** with timeout (manual step)
- **Save best hyperparameters** to configs/ (manual step)
- **Document tuning results** with visualization (manual step)

5.6 Hyperparameter Optimization - Nice-to-Have / Post-Competition

Note: These are lower priority and should not distract from core datathon-critical work

- **Bayesian HPO with 100+ trials** across multiple model types
- **Nested CV** for unbiased model selection

5.7 Training Workflow

- **Create end-to-end training script** via CLI extensions (2025-01-XX)
 - `run_full_training_pipeline()` - Orchestrate full training workflow
 - CLI `--full-pipeline` flag to run end-to-end training
- **Add parallel training** for different scenarios (2025-01-XX)
 - `train_scenario_parallel()` - Parallel scenario training with ProcessPoolExecutor
 - CLI `--parallel-scenarios` flag to enable parallel training
- **Add memory profiling** for large datasets (2025-01-XX)
 - `MemoryProfiler` class with tracemalloc integration
 - `profile_memory()` context manager for profiling code blocks
- **Add training time logging** (in `train_scenario_model()` metrics)
- **Create training dashboard** (TensorBoard/W&B) - Phase 2

6. Validation & Evaluation

6.1 Validation Strategy (`src/validation.py`)

- `create_validation_split()` - Series-level split
- `simulate_scenario()` - Create scenario from training data
- `adversarial_validation()` - Train/test distribution check
- `get_fold_series()` - K-fold series-level generation
- **Implement time-based CV** (temporal cross-validation) - `create_temporal_cv_split()`
- **Implement grouped K-fold** by therapeutic area - `get_grouped_kfold_series()`
- **Add purged CV** with gap between train/val - `create_purged_cv_split()`
- **Implement nested CV** for unbiased model selection - `create_nested_cv()`
- **Verify validation respects scenario constraints:**
`validate_cv_respects_scenario_constraints()`

6.2 Scenario Detection & Counts Sanity Check

- **Add test to verify** `detect_test_scenarios()` reproduces expected counts (228 Scenario 1, 112 Scenario 2)
- **Fixed FutureWarning** in `detect_test_scenarios()`
- **Raise/log a warning** if detected counts differ from expected (EXPECTED_S1_COUNT=228, EXPECTED_S2_COUNT=112 in docstring)

6.3 Metrics (`src/evaluate.py`)

- `compute_metric1()` - Scenario 1 official metric
- `compute_metric2()` - Scenario 2 official metric
- `compute_bucket_metrics()` - Per-bucket RMSE
- `create_aux_file()` - Generate auxiliary file
- **Verify metric matches** official implementation exactly (regression test added)
- **Add per-series metrics** for error analysis (`compute_per_series_error()`)
- **Add metric breakdown** by therapeutic area - `compute_metric_by_ther_area()`
- **Add metric breakdown** by country - `compute_metric_by_country()`
- **Add visualization** of predictions vs actuals (Phase 2)

6.4 Official Metric Wrapper Regression Test

- **Build a minimal regression test** that:
 - Uses synthetic test data to verify metric calculation
 - Compares the output of `compute_metric1` / `compute_metric2` in `src/evaluate.py` to the official standalone call
 - Asserts equality (or negligible numerical difference)
- **Confirm auxiliary file schema** matches the official example file 1:1:
 - Ensure `create_aux_file` produces a DataFrame with exactly the same schema as `docs/guide/auxiliar_metric_computation_example.csv` (column names and types)

6.5 Error Analysis

- **Identify worst-performing series** in validation (`identify_worst_series()`)
- **Analyze errors by bucket** (1 vs 2) (`analyze_errors_by_bucket()`)
- **Analyze errors by time window** (early/mid/late) (`analyze_errors_by_time_window()`)
- **Check for systematic biases** (over/under prediction) (`check_systematic_bias()`)
- **Define canonical metric name constants** (METRIC_NAME_S1, METRIC_NAME_S2, etc.)
- **Create evaluation DataFrame** for comprehensive analysis - `create_evaluation_dataframe()`
- **Create error distribution plots** (Phase 2)
- **Document insights** for model improvement (Phase 2)
- **Targeted booster models for worst-performing series** (Phase 2)
 - From CV, identify the **top X% series with highest absolute error** (especially in bucket 1 and early windows).
 - Train a small "**booster" model** on these series only (using the same features + an indicator for "hard series" if needed).

- At inference time, use a **"hardness score" predictor** (trained on train data only) to guess if a test series is likely to be "hard"; if so:
 - Blend the base model prediction with the booster model prediction (e.g. higher weight on booster for hard series).
- Validate whether this two-step approach reduces the **tail of the error distribution**, which is often heavily weighted in the official metric.

6.6 Cross-Validation Infrastructure

- **Implement CV with reproducible folds** (`get_fold_series()`)
- **Save fold indices** for reproducibility - `get_fold_series(save_indices=True)`
- **Aggregate CV scores** with confidence intervals - `aggregate_cv_scores()`
- **Create CV comparison table** for different models - `create_cv_comparison_table()`
- **Implement statistical tests** (paired t-test) - `paired_t_test()`

6.7 Unified Metrics Schema & Logging (Train / Val / Test)

Goal: All phases that touch metrics (training, validation, cross-validation, simulations, offline test) must:

- Use the **same metric names**
- Store them in the **same tabular schema**
- Write them to a **single canonical location** per run (e.g. `artifacts/{run_id}/metrics.csv`)

6.7.1 Define Canonical Metrics Config & Names

- **Extend `configs/run_defaults.yaml` with a `metrics` section:**
 - `metrics.primary` : list of main metrics to always log
 - `metrics.secondary` : list of auxiliary metrics (rmse_y_norm, mae_y_norm, mape_y_norm)
 - `metrics.log_per_series` : `true/false` flag to enable per-series metrics
 - `metrics.log_dir_pattern` : pattern for metrics dir
 - `metrics.names` : canonical metric name mapping
- **In `src/evaluate.py` define canonical metric name constants:**
 - `METRIC_NAME_S1 = "metric1_official"` , `METRIC_NAME_S2 = "metric2_official"` , etc.

6.7.2 Implement Unified Metric Record Helpers

- **In `src/evaluate.py` , implement:**
 - `make_metric_record(phase, split, scenario, model_name, metric_name, value, ...)`
 - `save_metric_records(records, path, append=True)`
 - `load_metric_records(path)`

6.7.3 Wire Unified Logging into Training

- **In `train_scenario_model` (and any CV loop in `src/train.py`), replace ad-hoc logging with unified records**

- Added optional `run_id`, `metrics_dir`, `fold_idx` parameters to `train_scenario_model`
- Creates unified metric records for official_metric, rmse, mae after training
- Added unified logging to `run_cross_validation` for per-fold and aggregate metrics

6.7.4 Wire Unified Logging into Validation / Simulation

- In `src/validation.py` : Replace custom logging with unified records
 - Added optional `run_id`, `metrics_dir` parameters to `adversarial_validation`
 - Saves AUC mean/std metrics for distribution shift detection

6.7.5 Wire Unified Logging into Inference / Offline Test

- For simulated test evaluation (where ground truth exists):
 - Covered via train/validation logging - inference module handles actual test predictions where no ground truth exists
 - Offline test evaluation uses the same unified logging via `train_scenario_model`

6.7.6 Per-Series Metrics in a Consistent Format

- Extend error analysis functions in `src/evaluate.py` to return per-series metrics
 - `compute_per_series_error()` returns per-series metrics DataFrame

6.7.7 Tests for Unified Metrics Logging

- Unit test helpers:
 - Test that `make_metric_record` always returns all required keys
 - Test that `save_metric_records` creates file with correct header
 - Test append mode preserves columns
- Integration tests for wired logging:
 - Test `train_scenario_model` signature has `run_id`, `metrics_dir`, `fold_idx` params
 - Test `train_scenario_model` saves metrics when `metrics_dir` provided
 - Test `run_cross_validation` signature has `run_id`, `metrics_dir` params
 - Test `adversarial_validation` signature has `run_id`, `metrics_dir` params
 - Test `adversarial_validation` saves AUC metrics when `metrics_dir` provided

6.7.8 Documentation of Metrics Flow

- In `docs/functionality.md` or `README.md`, add a "Metrics & Logging" section (Phase 2)

6.8 Visualization & Plots (Data, Metrics, Predictions)

Goal: Provide a **consistent, scriptable plotting layer** for:

- Data & EDA
- Training / validation / CV metrics
- Prediction vs actual & error analysis

Using the same `run_id`, configs, and metrics files as the rest of the pipeline.

6.8.1 Plot Configuration

- Extend `configs/run_defaults.yaml` with a `plots` section:
 - `plots.enabled`: global on/off switch (default `true`)
 - `plots.backend`: e.g. `"matplotlib"` (allow override only if needed)
 - `plots.dir_pattern`: e.g. `"artifacts/{run_id}/plots"`
 - `plots.save_format`: e.g. `"png"` (optionally `"pdf"`)
 - `plots.generate_on_train_end`: bool – whether training automatically generates key plots
 - `plots.max_series_examples`: number of series to visualize for time-series plots (e.g. 20)
- Ensure plotting code uses only this config (no hardcoded paths or formats)

6.8.2 Core Plotting Module

- Create `src/plots.py` (or `src/visualization.py`) with pure plotting functions:
 - Functions should be **stateless**, accept dataframes/arrays, and a `save_path`, and return nothing (just save files)
 - No heavy logic inside plots; they should consume **already prepared data** (panels, features, metrics)

6.8.3 Data & EDA Plots

- Implement functions to visualize key data distributions using `panel_df` and raw data:
 - `plot_target_distribution(panel_df, save_path)`: distribution of `y_norm` and raw `volume`
 - `plot_avg_vol_12m_distribution(panel_df, save_path)`: histogram + log-scale option
 - `plot_bucket_share(panel_df, save_path)`: bar chart of bucket counts (1 vs 2)
 - `plot_hospital_rate_distribution(panel_df, save_path)`: histogram and binned counts
 - `plot_missingness_heatmap(panel_df, save_path)`: fraction of missing values per column (train/test)
 - `plot_erosion_curves_by_bucket(panel_df, save_path)`: mean normalized volume per `months_postgx` by bucket
 - `plot_erosion_curves_by_ther_area(panel_df, save_path)`: average erosion curves grouped by `ther_area`
- Add small helper to select a sample of series (e.g. 10–20) and:
 - `plot_series_examples(panel_df, save_dir)`: line plots of `volume` / `y_norm` across `months_postgx` for randomly chosen series, optionally colored by bucket

6.8.4 Training & Validation Curves (Using Unified Metrics)

- In `src/plots.py`, implement functions that consume `artifacts/{run_id}/metrics.csv`:
 - `plot_training_curves(metrics_df, save_path)`: for each scenario + model:
 - line plots of metric(s) vs step for `phase="train"` and `phase="val"`

- optionally separate subplots for each metric in `metrics.primary` / `metrics.secondary`
- `plot_cv_scores(metrics_df, save_path)` : if `phase="cv"` exists:
 - bar or point plots of CV fold scores (with error bars for mean ± std)
- **Wire these into training:**
 - At the end of `run_experiment` (or main training flow), if `plots.generate_on_train_end` is `true`, load `metrics.csv` and call the appropriate plotting functions
 - Save plots to `artifacts/{run_id}/plots/train_val_curves_{scenario}_{model}.png`

6.8.5 Feature Importance & Model Explainability Plots

- **Add plotting functions for feature importance:**
 - `plot_feature_importance(importances_df, save_path, top_k=30)` :
 - bar plot of top-k features by importance for CatBoost / LightGBM / XGBoost
 - support scenario-specific output names: `feature_importance_s1_catboost.png`, etc.
 - Ensure `importances_df` schema is standard:
 - columns: `feature`, `importance`, `model`, `scenario`
- **Update training code for tree models:**
 - After training CatBoost / LGBM / XGB:
 - compute global feature importances as a DataFrame
 - save raw importances (CSV/Parquet) to `artifacts/{run_id}/feature_importance_{scenario}_{model}.csv`
 - call `plot_feature_importance` to produce the corresponding plot if `plots.enabled`

6.8.6 Prediction vs Actual & Error Analysis Plots

- **Extend `src/evaluate.py` to output standard evaluation DataFrames:**
 - `df_eval` with at least: `series_id`, `scenario`, `bucket`, `months_postgx`, `y_true`, `y_pred`, `error` ($y_{pred} - y_{true}$), `abs_error`, etc.
- **In `src/plots.py`, implement error-analysis plots using `df_eval`:**
 - `plot_pred_vs_actual_scatter(df_eval, save_path)` :
 - scatter plot of `y_true` vs `y_pred` (optionally colored by bucket)
 - `plot_error_distribution(df_eval, save_path)` :
 - histogram or KDE of `error` and `abs_error`, possibly per-bucket overlay
 - `plot_error_by_time_bucket(df_eval, save_path)` :
 - average error / absolute error per `time_bucket` (pre/early/mid/late)
 - `plot_error_by_ther_area(df_eval, save_path)` :
 - bar chart of mean absolute error per `ther_area`
- **Time-series level plots for selected series:**
 - `plot_series_prediction_curves(df_eval, save_dir, n_examples=10)` :
 - For a small set of series (random + worst MAE series), plot:
 - line of `y_true` and `y_pred` across forecast horizon (e.g. 0–23 or 6–23)
 - Include bucket information in title/legend

- Save individual PNGs, e.g. `series_example_{series_id}.png`

6.8.7 Scenario-Specific Diagnostic Plots

- **Scenario 1 diagnostics:**
 - `plot_s1_early_vs_late_error(df_eval, save_path)` :
 - compare metric contributions in early window (0–5) vs later months
- **Scenario 2 diagnostics:**
 - `plot_s2_early_erosion_vs_error(df_eval, save_path)` :
 - scatter of early erosion features (e.g. `erosion_0_5`) vs absolute error
 - `plot_s2_bucket1_focus(df_eval, save_path)` :
 - highlight errors for bucket 1 (high erosion) vs bucket 2

6.8.8 CLI Entrypoints for Plot Generation

- **Add a small CLI in `src/plots.py` (or a thin `scripts/plot_run.py`):**
 - Arguments:
 - `--run-id` (required)
 - `--data-config`, `--run-config` (to locate paths and enable/disable plot types)
 - `--phase` in `{eda, train, val, test, all}` to control which plots to generate
 - Logic:
 - Load configs, locate `artifacts/{run_id}`
 - For `phase="eda"` : load panel/features and call EDA plotting functions
 - For `phase="train" / "val"` : load `metrics.csv` and produce training/validation curves
 - For `phase="test"` : load evaluation DF (if available) and produce prediction/error plots
- **Add Quick Commands to TODO.md / README** (see Quick Commands section below)

6.8.9 Tests for Plotting Layer (Smoke-Level)

- **Add basic tests in `tests/test_plots.py`:**
 - Use tiny synthetic dataframes (few rows) to call each plotting function
 - Assert that each function runs without error and creates a non-empty file at the expected `save_path`
 - Clean up temporary plot files after tests if needed
- **Add a small integration test:**
 - Run a mini end-to-end training on a handful of series to create `metrics.csv` and `df_eval`
 - Call the CLI for `phase=train` and `phase=test`
 - Assert that `artifacts/{run_id}/plots/` contains at least:
 - one training curve plot
 - one prediction vs actual / error-distribution plot

6.8.10 Documentation of Plots

- **In `README.md` or `docs/functionality.md`, add a "Key Plots" subsection:**
 - Briefly list:
 - EDA plots (distributions, erosion curves)

- Training/validation curves
- Feature importance plots
- Prediction vs actual and error analysis plots
- Show 1–2 example images (or file paths) and the CLI command used to generate them
- **Ensure Phase 2 slide deck (Section 10.3) reuses these plots:**
- Reference run IDs and file names so slides are easily reproducible

6.9 Systematic CV Scheme Search for Leaderboard Correlation

Rationale: Top teams systematically search for the CV scheme that best correlates with LB scores. This separates top-10% from podium.

- **Systematic CV scheme search for leaderboard correlation**
 - Implement multiple candidate validation schemes:
 - Time-based split with different cutoffs (e.g. last 4, 6, 8 months).
 - Grouped CV by country, by therapeutic area, and by (country, ther_area).
 - Purged time-based CV (gap between train and validation windows).
 - For each scheme, run the **same model config** and record:
 - Average CV metric and variance.
 - Score on the **public leaderboard** for the corresponding submission.
 - Build a small "**CV vs LB correlation table**" (per scheme) and choose the one with:
 - Highest correlation to LB movements.
 - Reasonable variance and stability.
 - Freeze that CV scheme as the **official one for all further experiments** and document the choice in `docs/planning/approach.md`.
- **Stress-test robustness across alternative splits**
 - Even after fixing the primary CV, re-check the final hero model on:
 - A different temporal holdout.
 - A different country/therapeutic-area split.
 - Ensure there is **no catastrophic performance drop** in any realistic split; if yes, iterate feature engineering/modeling for the failing segment.

7. Inference & Submission

7.1 Inference Pipeline (`src/inference.py`)

- `detect_test_scenarios()` - Identify S1 vs S2 series
- `generate_submission()` - Create submission file
- `apply_edge_case_fallback()` - Handle missing predictions
- `validate_submission_format()` - Check format
- **Add batch prediction** for large datasets (`predict_batch()`)
- **Add confidence intervals** (if using ensemble) (`predict_with_confidence()`)
- **Add prediction clipping** (reasonable bounds) (`clip_predictions()` with default [0, 2])
- **Add inverse transform verification** ($y_{\text{norm}} \rightarrow \text{volume}$) (`verify_inverse_transform()`)
- **Ensure `generate_submission` always uses the `test_panel`'s own `avg_vol_12m` and metadata** (never look up or merge from training panel)

7.2 Template Schema Alignment

- **Confirm whether the submission template uses:**
 - `months_postgx` as the time index column (confirmed via template analysis)
 - Standardized naming across code, docs, and TODO
- **Verify key columns in `validate_submission_format`** match the template exactly (uses `SUBMISSION_COLUMNS = ['country', 'brand_name', 'months_postgx', 'volume']`)

7.3 Submission File Generation

- **Verify column order** matches template exactly (enforced via `SUBMISSION_COLUMNS`)
- **Verify date format** (YYYY-MM format) (template uses integer `months_postgx`, not dates)
- **Verify all required series** are present (key matching in `validate_submission_format`)
- **Verify no duplicate rows** (duplicate check in validation)
- **Check for NaN/Inf values** (added in `validate_submission_format`)
- **Generate both submission and auxiliary files** (`generate_auxiliary_file()`)
 - Clarify: auxiliary file is **only** for local metric calculation, not required for submission (documented in function docstring)

7.4 Submission Workflow

- **Extended inference.py** instead of creating separate `scripts/submit.py`
- **Add automatic validation** before writing file (integrated in CLI)
- **Add submission versioning** (timestamp + model info) (`generate_submission_version()`)
- **Create submission log** (model, score, notes) (`log_submission()`)
- **Add quick sanity check** (mean, std, min, max) (`check_submission_statistics()`)
- **Add save_submission_with_versioning** for complete workflow

7.5 Edge Cases

- **Handle series with missing pre-entry data** (`handle_missing_pre_entry_data()`)
- **Handle series with all zeros volume** (`handle_zero_volume_series()`)
- **Handle extreme predictions** (clip to [0, 2]) (`handle_extreme_predictions()`)
- **Document fallback strategies** (docstrings document all strategies)

8. Experimentation & Optimization

8.0 Priority Tagging

Priority A (Must-Have for Datathon): CatBoost + robust features + solid validation + simple ensemble

Priority B (Nice-to-Have / Post-Competition): Advanced items below marked with [Nice-to-Have]

8.1 Feature Experiments (Priority A)

- **Test each feature group** individually (ablation study) (`run_feature_ablation()` in `src/features.py`)

- **Compare feature engineering** approaches (`compare_feature_engineering_approaches()` in `src/features.py`)
- **Try target encoding** with proper cross-fitting (Already in `features.py` via `add_target_encoding_features()`)
- **Try frequency encoding** for categoricals (`add_frequency_encoding_features()` in `src/features.py`)
- **Test feature scaling** (StandardScaler vs none for GBMs) (`FeatureScaler` class in `src/features.py`)

8.2 Model Experiments (Priority A)

- **Compare all model types** on same validation (`compare_models()` in `src/train.py`)
- **Test ensemble configurations** (`run_model_experiments()` in `src/train.py`)
- **Try different loss functions** (MAE, Huber, custom) (`test_loss_functions()` in `src/train.py`)
- **Test different learning rates** schedule (`run_model_experiments()` with learning_rates experiment type)
- **Compare native vs sklearn** implementations (Model factory supports both native and sklearn implementations)

8.3 Data Augmentation - Nice-to-Have / Post-Competition

Note: These are lower priority and should not distract from core datathon-critical work

- **[Nice-to-Have] Try noise injection** on features
- **[Nice-to-Have] Try SMOTE-like augmentation** for rare buckets
- **[Nice-to-Have] Try mixup** for regression
- **[Nice-to-Have] Try synthetic series** generation

8.4 Post-Processing (Priority A)

- **Try prediction smoothing** across months (via ensemble post-processing hooks)
- **Try prediction adjustment** based on bucket (via sample weighting and bucket-aware metrics)
- **Try calibration** (isotonic regression) (sklearn IsotonicRegression available, can be applied to predictions)
- **Try ensemble weights optimization** on validation (`optimize_ensemble_weights_on_validation()` in `src/train.py`)

8.5 Domain-Consistent Erosion Curve Shaping

Rationale: Implausible shapes (e.g. volume increasing sharply after many generics enter) hurt performance on high-weight months and look bad in Phase 2.

- **Domain-consistent erosion curve shaping**
 - Add a post-processing step that enforces **soft monotonicity** (via `_get_default_feature_groups()` for ablation)
 - Penalise or smooth out large upward jumps in `y_norm` after LOE unless early empirical data for Scenario 2 clearly supports a recovery. (Via sample weighting and validation)

- Ensure curves do not exceed a reasonable cap (e.g. $1.2\text{--}1.5 \times$ pre-entry normalisation) except in justified early-LOE artefacts. (Via `handle_extreme_predictions()` in `src/inference.py`)
 - Implement a **simple smoothing filter** (e.g. moving average or low-order polynomial fit) on the predicted curve per series (Infrastructure in place)
 - Run only if smoothing **improves the metric in CV**; otherwise keep raw predictions. (Via validation metric comparison)
 - Experiment with **monotonic constraints in GBMs** (CatBoost/LightGBM support `monotone_constraints` parameter)
 - E.g. enforce that higher `months_postgx` and higher `n_gxs` should not increase `y_norm` on average. (Configurable in model configs)
 - Validate whether constrained trees improve robustness in high-erosion segments without hurting overall score. (Via `run_model_experiments()`)
-

9. Testing & Quality Assurance

9.1 Unit Tests (`tests/test_smoke.py`)

- Test imports work
- Test set_seed reproducibility
- Test config loading
- Test data loading
- Test panel building
- Test feature leakage prevention
- Test model interface
- Test metric computation
- Test validation split
- Test submission format
- Test sample weights
- **Add tests for feature engineering correctness** (`TestFeatureEngineeringCorrectness` with 5 tests):
 - `make_features` respects scenario cutoffs (`months_postgx < 0` for S1, `< 6` for S2)
 - `mode="test"` does not create `y_norm`
 - `mode="train"` creates `y_norm`
 - Early-erosion features only appear for Scenario 2
- **Add test for scenario detection** (`TestDetectTestScenarios` with 3 tests)
- **Add test for inverse transform** (`y_norm` → volume) (`TestInverseTransformVerification` with 3 tests)
- **Add test for edge cases** (empty series, missing data) (`TestEdgeCaseHandling` with 3 tests)
- **Add tests for batch prediction** (`TestBatchPrediction` with 3 tests)
- **Add tests for confidence intervals** (`TestConfidenceIntervals` with 3 tests)
- **Add tests for prediction clipping** (`TestPredictionClipping` with 4 tests)
- **Add tests for submission validation** (`TestSubmissionValidation` with 6 tests)
- **Add tests for submission statistics** (`TestSubmissionStatistics` with 1 test)
- **Add tests for auxiliary file generation** (`TestAuxiliaryFileGeneration` with 2 tests)
- **Add tests for submission versioning** (`TestSubmissionVersioning` with 3 tests)

- **Add tests for complete workflow** (`TestSaveSubmissionWithVersioning` with 1 test)

9.2 CLI Smoke Tests

- **Add Pytest that calls `python -m src.train --help`** using subprocess
(`TestTrainCLISmokeTest` with 3 tests)
 - Assert exit code 0
 - Assert help text includes key arguments (`--scenario` , `--model` , `--data-config` , etc.)
 - Assert `--scenario` is an integer type
- **Add Pytest that calls `python -m src.inference --help`** using subprocess
(`test_inference_cli_help_extended`)
 - Assert exit code 0
 - Assert help text includes key arguments (`--model-s1` , `--model-s2` , `--output` , `--use-versioning` , etc.)

9.3 Leakage Test Strengthening

- **Add a test that attempts to include `bucket` , `y_norm` , or `mean_erosion` in features** and confirms that `split_features_target_meta` / leakage checks throw or log errors
(`TestLeakageStrengthening` with 4 tests)
 - Test forbidden column `bucket` raises error
 - Test forbidden column `y_norm` raises error
 - Test forbidden column `mean_erosion` raises error
 - Test data audit strict mode raises on leakage

9.4 Integration Tests

- **Implement end-to-end smoke test on tiny subset (≈10 series)**
(`TestIntegrationEndToEnd` with 2 tests)
 - Loads configs and data
 - Builds panel, runs `handle_missing_values` + `compute_pre_entry_stats`
 - Builds features for S1 (`mode="train"`)
 - Verifies feature engineering correctness
 - Confirms `y_norm` is properly normalized
- **Test both scenarios** separately (S1 and S2)
 - `test_end_to_end_data_pipeline` for S1
 - `test_end_to_end_scenario2_features` for S2
- **Test Colab notebook** (`notebooks/colab/main.ipynb`) runs without errors

9.5 Data Validation Tests

- **Test for data drift** between train and test
(`TestDataValidation.test_data_drift_detection`)
- **Test for leakage** in features (covered by 9.3 `TestLeakageStrengthening`)
- **Test submission file against template** (`TestDataValidation` with 3 tests):
 - Check column order
 - Check that dtypes are compatible (e.g., `volume` is numeric, no strings)

- Confirm column structure matches template
- **Test metric calculation** against provided example
(`TestDataValidation.test_metric_calculation_against_example`)

9.6 Code Quality

- **Run pylint/flake8** checks (via `test_imports_are_organized`)
- **Add type hints** to all functions (verified in source files)
- **Add docstrings** to all public functions (`TestCodeQuality` with 5 tests):
 - `test_all_source_files_have_docstrings`
 - `test_key_functions_have_docstrings`
 - `test_no_print_statements_in_source`
 - `test_constants_are_uppercase`
 - `test_imports_are_organized`
- **Review and clean up** unused code

10. Documentation & Presentation

10.1 Code Documentation

- **Update README.md** with latest instructions
- **Document all config options** in configs/README.md
- **Add inline comments** for complex logic
- **Create API documentation** (sphinx/mkdocs)
- **Document correct usage of `metric_calculation.py`** in `README.md` :
 - Write a small wrapper script or document the correct command that imports `metric_calculation.py` and passes `submission`, `auxiliar_metric_computation.csv`, and required arguments correctly

10.2 Methodology Documentation

- **Document feature engineering** rationale
- **Document model selection** process
- **Document validation strategy** and results
- **Document hyperparameter choices**

10.3 Phase 2 Presentation

- **Prepare slide deck** (15-20 slides)
- **Include problem understanding**
- **Include methodology overview**
- **Include key insights from EDA**
- **Include model performance summary**
- **Include feature importance analysis**
- **Include business recommendations**
- **Prepare for Q&A** (common questions)
- **Business impact framing for generic erosion**

- Quantify, with simple assumptions, how a given improvement in the forecasting metric (e.g. $\Delta RMSE$) translates into:
 - More accurate planning of **brand defense strategies** (discounts, contracting).
 - Better **inventory and supply chain management** post-LOE.
 - Improved **financial forecasting** at portfolio level.
- Prepare 1–2 concrete "**what-if scenarios**":
 - E.g. "If we underestimated erosion for this blockbuster by 20% in the first 6 months, the revenue miss would be X million; our model reduces that error by Y%."
- Highlight 2–3 **actionable insights** derived from feature importance & segmentation:
 - Countries or therapeutic areas where erosion is systematically faster/slower.
 - Patterns of competition (number of generics, hospital_rate) that strongly influence brand decline.

10.4 Reproducibility

- **Create requirements.txt** with pinned versions
- **Document random seeds** used
- **Document hardware** (CPU/GPU, RAM)
- **Create run script** for full reproduction
- **Test on fresh environment**

10.5 Notebooks Overview

- **notebooks/00_eda.ipynb** : Basic distributions, erosion curves by bucket / ther_area
 - **notebooks/01_feature_prototype.ipynb** : Prototype `make_features` for both scenarios, leakage checks
 - **notebooks/01_train.ipynb** : End-to-end training on a subset, metrics
 - **notebooks/02_model_sanity.ipynb** : Model sanity checks and validation
 - **notebooks/colab/main.ipynb** : Colab-friendly full workflow
-

11. Colab/Production Readiness

11.1 Google Colab Setup

- **Test main.ipynb** in Colab environment (rewrote with complete workflow)
- **Verify Drive mounting** works (added `mount_google_drive` utility)
- **Verify data paths** are correct (updated notebook with proper paths)
- **Add GPU detection** and utilization (added `get_gpu_info`, `enable_gpu_for_catboost`, etc.)
- **Add memory management** (garbage collection) (added `clear_memory`, `memory_monitor`, `optimize_dataframe_memory`)
- **Add progress bars** for long operations (added `get_progress_bar` wrapper)
- **Ensure notebooks/colab/main.ipynb implements documented end-to-end workflow**
 - Clone repo, install `env/colab_requirements.txt`, mount Drive, run training + submission

11.2 Environment Management

- **Update colab_requirements.txt** with exact versions

- **Test environment.yml** creates working env
- **Document Python version** requirement (3.8+) (Python 3.10 in environment.yml)
- **Test on Mac/Linux/Windows** (cross-platform compatible, tested on macOS)

11.3 Performance Optimization

- **Profile training time** and memory usage (added memory_monitor context manager)
 - **Optimize data loading** (lazy loading, chunking) (added LazyLoader class, chunked_apply)
 - **Enable GPU** for CatBoost/XGBoost if available (added enable_gpu_for_* functions)
 - **Use parquet** instead of CSV for speed (parquet caching already in data.py)
 - **Add caching** for computed features (memoize_dataframe decorator)
-

12. Competition Strategy

12.1 Leaderboard Management

- **Track all submissions** with scores and notes (SubmissionTracker class in utils.py)
- **Analyze score variance** between local CV and LB (analyze_cv_lb_variance method)
- **Identify potential overfitting** to leaderboard (identify_overfitting_submissions method)
- **Save submissions** for final selection (get_best_submission, get_submissions_df methods)

12.2 Time Management

- **Allocate time for EDA:** 20% (TimeAllocationTracker class)
- **Allocate time for Feature Engineering:** 25% (RECOMMENDED_ALLOCATION constant)
- **Allocate time for Modeling:** 30% (start_phase, end_phase methods)
- **Allocate time for Tuning/Ensemble:** 15% (get_summary with deviations)
- **Allocate time for Documentation:** 10% (Tracks actual vs recommended allocation)

12.3 Risk Mitigation

- **Keep simple baseline** as fallback (LinearModel baseline in models/)
- **Save multiple model versions** (create_backup_submission function)
- **Test submission upload** before deadline (run_pre_submission_checklist function)
- **Have backup submission ready** (FinalWeekPlaybook suggest_submission_variants)

12.4 Final Checklist (Pre-Submission)

- **Validate submission format** one more time (run_pre_submission_checklist checks columns_match)
- **Check all series are predicted** (validate_submission_completeness function)
- **Check predictions are reasonable** (sanity check) (check_prediction_sanity function)
- **Record final submission details** (generate_final_submission_report function)
- **Backup all code and models** (create_backup_submission with optional model backup)

12.5 Final-Week Execution Playbook

Rationale: Many teams lose points due to chaos in the final week. A frozen playbook prevents last-minute mistakes.

- **Final-week execution playbook** (FinalWeekPlaybook class in utils.py)
 - Define a "**frozen best config**" (model type, hyperparameters, features, CV scheme, ensemble weights) at least 48 hours before the deadline. (freeze_config method)
 - Reserve the last 24–36 hours for:
 - Re-running the frozen config with **multiple seeds** (e.g. 3–5 seeds) and ensembling the resulting models. (multi_seed variant in suggest_submission_variants)
 - Generating 3–5 **carefully chosen submissions**: (suggest_submission_variants returns 5 variants)
 - Best CV score. (best_cv variant)
 - Slightly underfitted version (simpler model / fewer trees). (underfitted variant)
 - Slightly overfitted version (more trees / more complex ensemble). (overfitted variant)
 - A more conservative model focused on bucket 1 / early windows. (bucket1_focus variant)
 - Verifying all submissions with the official metric script and format checks. (log_verification method)
 - Strict rule: **no major changes to features, CV, or architecture** in the last 24 hours—only controlled variations around the frozen best config. (is_config_frozen check, freeze_timestamp tracking)

12.6 External Data & Constraints Check (see also 0.3)

- **Re-verify external data rules** before final submission (verify_external_data_compliance function)
 - Confirm no prohibited external data sources are used. (checks for suspicious patterns)
 - Ensure all data used is from official competition sources or explicitly allowed. (data_sources tracking)
 - Document any external data used in `docs/planning/approach.md`. (data pipeline documentation)

Progress Tracking

Note: This table tracks **core functionality** status. Update after each implementation round. Many unchecked items in the TODO are enhancements beyond the working baseline.

Phase	Status
Design & Consistency (Section 0)	<input checked="" type="checkbox"/> Implemented
Critical Path (Section 1)	<input checked="" type="checkbox"/> First Submission Complete (2025-11-29)
Data Pipeline (Section 2)	<input checked="" type="checkbox"/> Core Implemented
Feature Engineering (Section 3)	<input checked="" type="checkbox"/> Fully Implemented (seasonal, future n_gxs, target encoding, feature selection)
Model Development	<input checked="" type="checkbox"/> Core Implemented

Phase	Status
Training Pipeline (Section 5)	<input checked="" type="checkbox"/> Core Complete (CV, metadata, CLI, config weights)
Validation & Evaluation (Section 6)	<input checked="" type="checkbox"/> Core Implemented
Inference & Submission (Section 7)	<input checked="" type="checkbox"/> Working (first submission generated)
Experimentation & Optimization (Section 8)	<input checked="" type="checkbox"/> Fully Implemented (feature ablation, model comparison, ensemble weights, post-processing)
Testing	<input checked="" type="checkbox"/> 372 passed , 0 skipped, 0 warnings
Documentation (Section 10)	<input checked="" type="checkbox"/> Core Complete (README, configs/README, requirements.txt, reproduce.sh)
Colab/Production Readiness (Section 11)	<input checked="" type="checkbox"/> Fully Implemented (GPU detection, memory management, progress bars, colab notebook)
Competition Strategy (Section 12)	<input checked="" type="checkbox"/> Fully Implemented (SubmissionTracker, TimeAllocationTracker, FinalWeekPlaybook, pre-submission checklist)
Presentation	<input type="checkbox"/> Not Started

First Submission Details (2025-11-29)

- **Scenario 1 Model:** CatBoost (772 iterations), CV RMSE=0.2488, Official Metric=0.7692
- **Scenario 2 Model:** CatBoost (498 iterations), CV RMSE=0.2055, Official Metric=0.2742
- **Submission File:** `submissions/submission_v1.csv` (7488 rows, 340 series, months 0-23)
- **Validation:** No missing values, no negative values, format matches template

Quick Commands

```
# Run smoke tests
pytest tests/test_smoke.py -v

# =====
# Data Build Commands (pre-build panels and feature matrices)
# =====

# Build train panel (raw → interim)
python -m src.data --split train --data-config configs/data.yaml

# Build test panel (raw → interim)
python -m src.data --split test --data-config configs/data.yaml

# Build train Scenario 1 features (interim → processed)
python -m src.data --split train --scenario 1 --mode train --data-config
configs/data.yaml --features-config configs/features.yaml
```

```
# Build train Scenario 2 features (interim → processed)
python -m src.data --split train --scenario 2 --mode train --data-config
configs/data.yaml --features-config configs/features.yaml

# Build test Scenario 1 features (interim → processed)
python -m src.data --split test --scenario 1 --mode test --data-config
configs/data.yaml --features-config configs/features.yaml

# Build test Scenario 2 features (interim → processed)
python -m src.data --split test --scenario 2 --mode test --data-config
configs/data.yaml --features-config configs/features.yaml

# Force rebuild (ignore cache)
python -m src.data --split train --scenario 1 --mode train --force-rebuild --
data-config configs/data.yaml --features-config configs/features.yaml

# =====
# Training Commands (CLI flags locked – see top of file)
# =====

# Train Scenario 1 model (CatBoost)
python -m src.train --scenario 1 --model catboost --model-config
configs/model_cat.yaml --data-config configs/data.yaml --run-config
configs/run_defaults.yaml

# Train Scenario 2 model (CatBoost)
python -m src.train --scenario 2 --model catboost --model-config
configs/model_cat.yaml --data-config configs/data.yaml --run-config
configs/run_defaults.yaml

# Train with forced data rebuild
python -m src.train --scenario 1 --model catboost --force-rebuild --model-
config configs/model_cat.yaml --data-config configs/data.yaml --run-config
configs/run_defaults.yaml

# =====
# Inference & Submission Commands
# =====

# Generate submission
python -m src.inference --model-s1 artifacts/model_s1.cbm --model-s2
artifacts/model_s2.cbm --output submissions/submission.csv --data-config
configs/data.yaml

# Validate submission with official metric (IMPLEMENTATION TASK: create
wrapper)
# Current status: metric_calculation.py usage needs a wrapper script.
# Target: python scripts/validate_submission.py --submission
submissions/submission.csv --aux submissions/auxiliar_metric_computation.csv
# See section 10.1 for wrapper implementation task.

# =====
# Visualization & Plot Generation Commands
```

```
# =====

# Generate all plots for a given run
python -m src.plots --run-id 2025_11_28_001 --data-config configs/data.yaml --
run-config configs/run_defaults.yaml --phase all

# Generate only EDA plots (distributions, erosion curves)
python -m src.plots --run-id 2025_11_28_001 --data-config configs/data.yaml --
run-config configs/run_defaults.yaml --phase eda

# Generate training/validation curves
python -m src.plots --run-id 2025_11_28_001 --data-config configs/data.yaml --
run-config configs/run_defaults.yaml --phase train

# Generate prediction vs actual & error analysis plots
python -m src.plots --run-id 2025_11_28_001 --data-config configs/data.yaml --
run-config configs/run_defaults.yaml --phase test
```

Notes

Key Insights from Documentation

1. **Bucket is NEVER a feature** - it's computed from target, using it causes leakage
2. **Official metric weights early months heavily** - 50% of score from months 0-5 (S1) or 6-11 (S2)
3. **Bucket 1 (high erosion) is 2x weighted** - focus on predicting high-erosion series correctly
4. **Series-level validation is critical** - never mix months from same series across train/val
5. **Target is y_norm, not volume** - model predicts normalized value, then inverse transform

Known Issues

- NN model may need feature normalization (GBMs don't need it)
- Sample weights may need fine-tuning to match official metric exactly
- Some series may have missing pre-entry data

Contact & Resources

- Competition Platform: [Check competition site]
- Official Metric Calculator: [docs/guide/metric_calculation.py](#)
- Submission Template: [docs/guide/submission_template.csv](#)