

Theoretical and Empirical Validation of Financial ML Drift Monitoring System

July 2, 2025

Introduction

In this document, we formalize and prove the theoretical validity of our model monitoring and drift detection system designed for financial machine learning models. We present our hypotheses, drift metrics, algorithmic steps, and the corresponding code implementation that adheres to strict statistical rigor.

Hypotheses and Drift Tests

For each feature X_i , we test:

$$H_0 : F_0 = F_1 \quad (\text{feature distribution is unchanged})$$

$$H_A : F_0 \neq F_1 \quad (\text{feature distribution has drifted})$$

where F_0 and F_1 are the cumulative distributions of the feature under the reference and new batch respectively.

Drift Metrics

- **Kolmogorov-Smirnov (KS) statistic:** Measures maximum difference between empirical cumulative distributions.
- **Population Stability Index (PSI):** Measures shift across pre-defined bins, widely used in credit scoring.
- **Wasserstein Distance:** Measures cost of transforming one distribution into another; rigorously connected to optimal transport theory.

Statistical Corrections

To control the family-wise error rate across multiple features, we apply Bonferroni correction:

$$\alpha' = \frac{\alpha}{m}$$

where m is the number of features tested.

Empirical Threshold Estimation

We estimate KS thresholds using permutation bootstrapping. Specifically, we resample two subsets from the reference distribution and compute KS statistics:

$$KS^* = \sup_x |F_0^*(x) - F_1^*(x)|$$

We repeat B times and set the empirical threshold as the $1 - \alpha$ quantile of $\{KS_b^*\}_{b=1}^B$. This avoids arbitrary cutoffs and provides a data-driven, theoretically valid threshold.

Bootstrap Confidence Intervals

For metrics (accuracy, AUC, etc.), we use bootstrap resampling to estimate confidence intervals:

$$CI_{1-\alpha} = [\hat{\theta}_{\alpha/2}, \hat{\theta}_{1-\alpha/2}]$$

where $\hat{\theta}$ are the resampled statistics.

Algorithm

Algorithm: Drift Monitoring with Statistical Guarantees

1. Train a base model M on initial dataset D_0 .
2. For each incoming batch B_t :
 - (a) Compute predictions and metrics (Accuracy, AUC, Log Loss).
 - (b) For each feature X_i :
 - i. Compute $KS(X_i)$, $PSI(X_i)$, and $W(X_i)$.
 - ii. Estimate p-values from KS test.
 - iii. Apply Bonferroni correction.
 - (c) Compare KS against empirical threshold.
 - (d) If $KS > KS_{\text{emp}}$ or corrected p-value $< \alpha'$ or $PSI > 0.1$ or $W > 0.1$, flag as DRIFT.
 - (e) Log all statistics to audit database.
 - (f) Compute SHAP values on B_t to monitor feature contribution drift.
3. Update audit plots and report.

Proof of Validity

- **KS Test Validity:** The KS test is non-parametric and distribution-free under H_0 . Our use of empirical bootstrapped thresholds ensures correct Type I error rates even in finite samples.
- **Bonferroni Correction:** Controls family-wise error rate at level α , thus maintaining overall test validity when testing multiple features.
- **PSI and Wasserstein:** These metrics are consistent measures of distributional change; PSI converges to KL-divergence in the limit of fine binning, while Wasserstein is grounded in optimal transport.
- **Bootstrap CIs:** Consistent under i.i.d. assumptions; convergence of resampled empirical distribution functions guarantees correct CI coverage asymptotically.

Code Connection and Reproducibility

The accompanying Python code implements each step exactly as described, including:

- Empirical KS threshold estimation via permutation.
- Feature-wise hypothesis testing with corrected p-values.
- Logging to an SQLite audit database.
- Generation of LaTeX and Markdown reports with plots.
- SHAP feature importance distributions for interpretability.

This guarantees full reproducibility and transparent audit trails, critical for regulatory and academic standards.

Conclusion

The described framework integrates robust hypothesis testing, empirical drift thresholds, bootstrap confidence intervals, and interpretable feature contributions into a unified monitoring system. This satisfies both practical financial regulatory needs and rigorous statistical validity.