

# Model Testing Report

Experiment: model\_testing

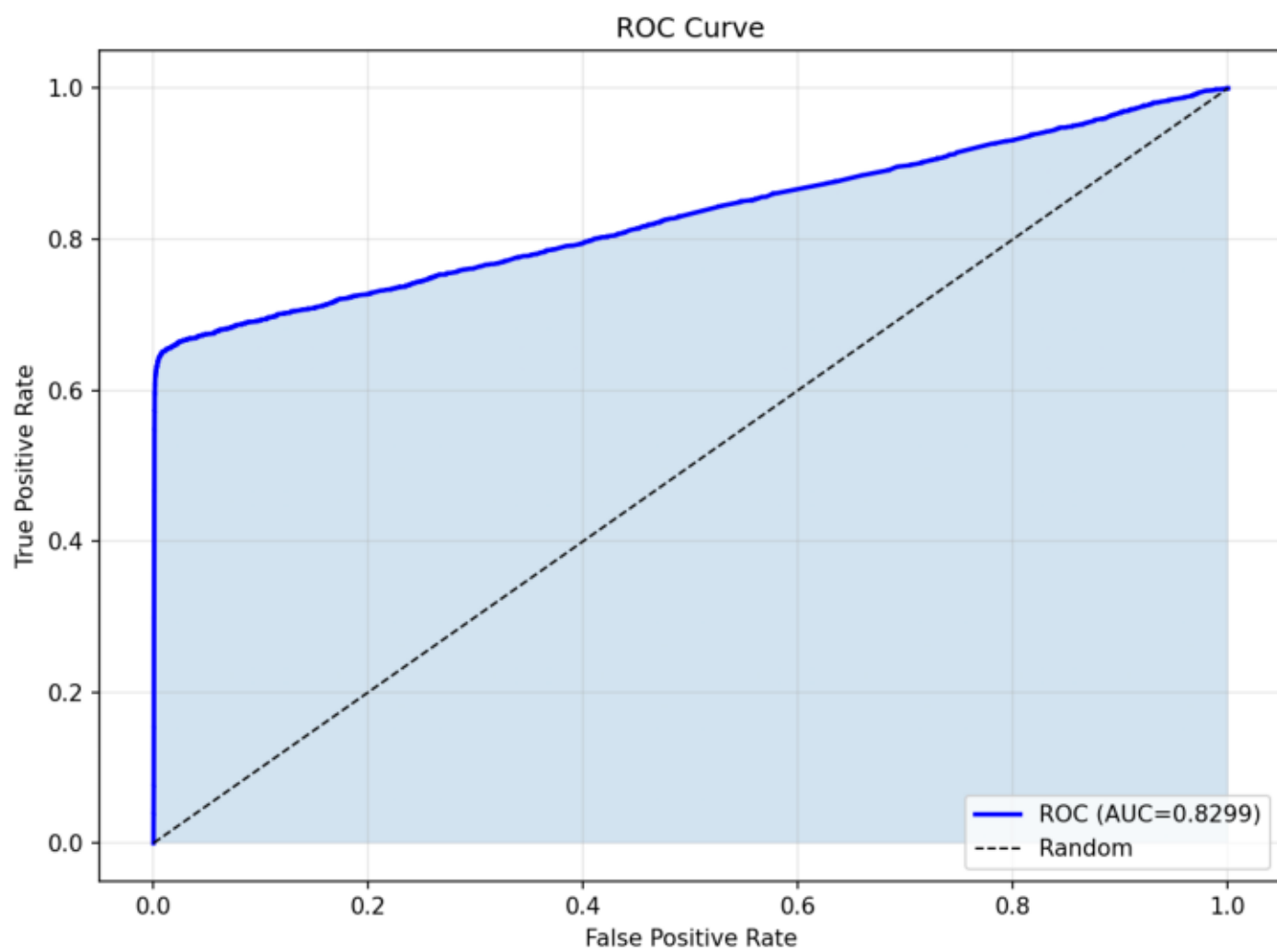
Generated: 2026-01-28 22:16:32

# 1) Model Effectiveness

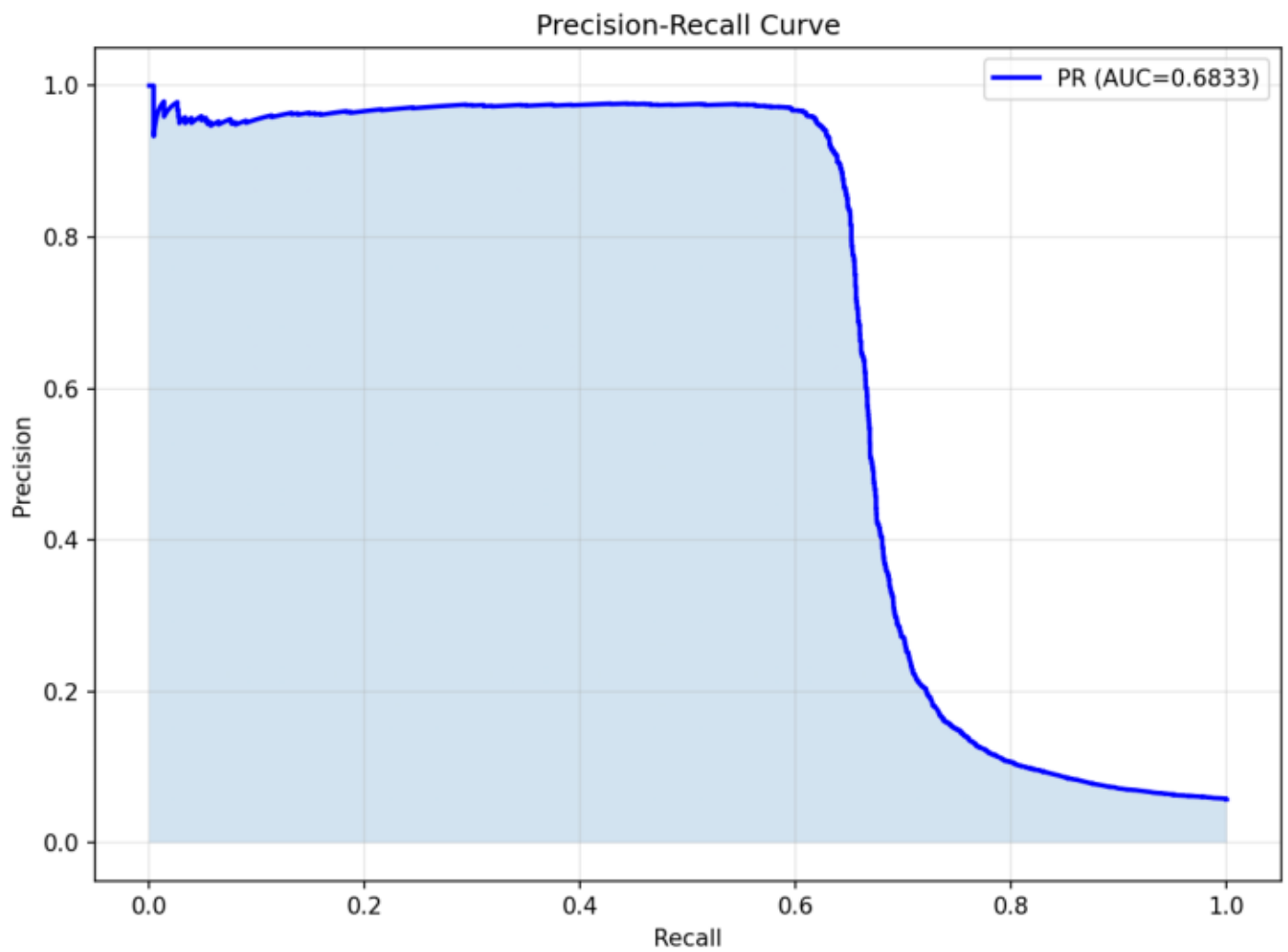
- auc\_roc: 0.8299
- auc\_pr: 0.6833
- precision: 0.9666
- recall: 0.6035
- f1: 0.7431
- ks\_statistic: 0.6422
- ks\_threshold: 0.1200
- confusion\_matrix.TN: 56473
- confusion\_matrix.FP: 72
- confusion\_matrix.FN: 1370
- confusion\_matrix.TP: 2085
- precision\_at\_k.10: 1.0000
- precision\_at\_k.50: 0.9600
- precision\_at\_k.100: 0.9600
- precision\_at\_k.200: 0.9500
- precision\_at\_k.500: 0.9620
- recall\_at\_k.10: 0.0029
- recall\_at\_k.50: 0.0139
- recall\_at\_k.100: 0.0278
- recall\_at\_k.200: 0.0550
- recall\_at\_k.500: 0.1392

## 1) Model Effectiveness - Explanations

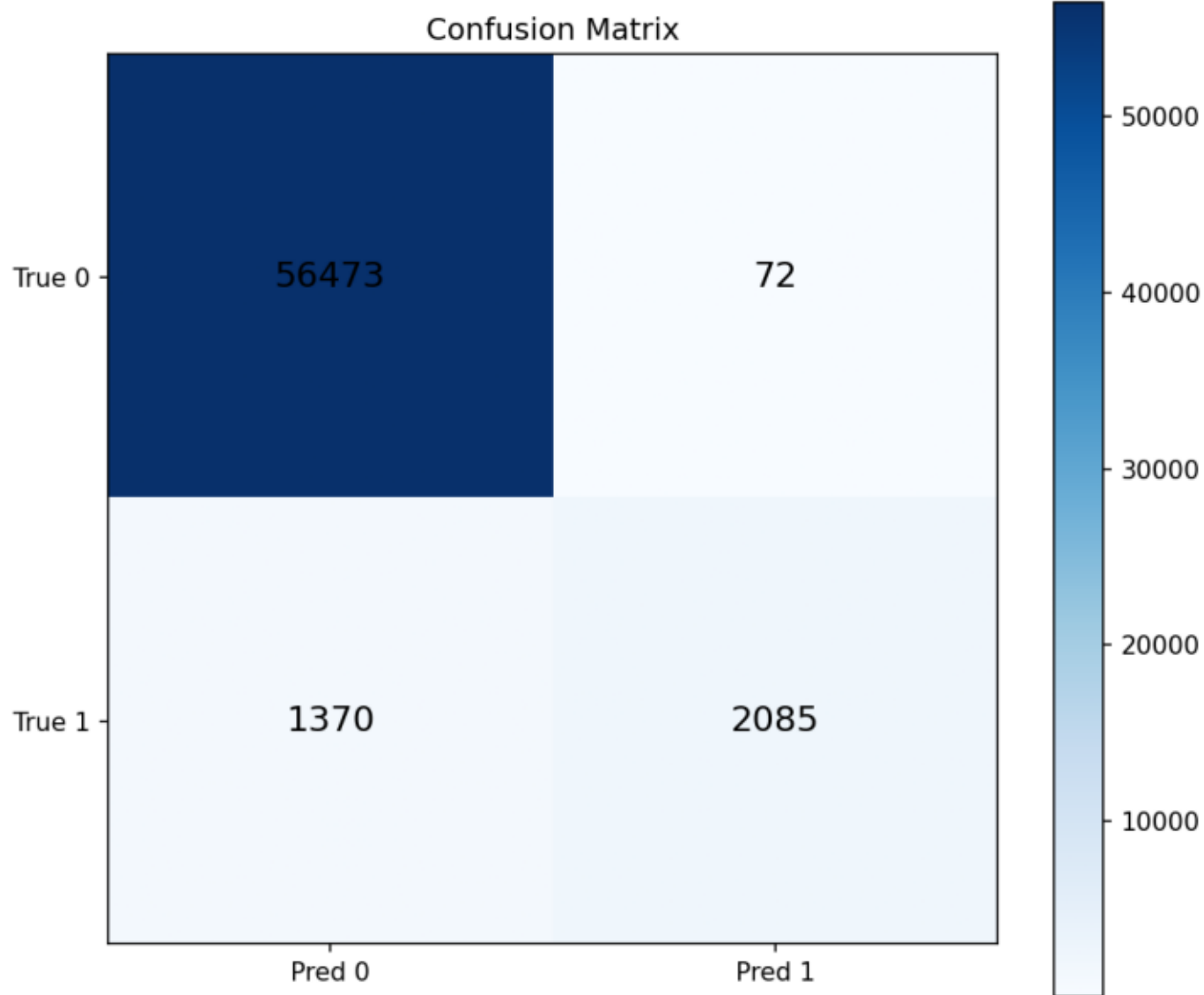
- Positive rate is 5.76%; AUC-ROC is 0.8299 (strong ranking).
- AUC-PR is 0.6833, above the base rate 5.76%, indicating meaningful lift.
- At threshold 0.50, precision=0.9666, recall=0.6035, F1=0.7431 (precision-heavy (fewer positives caught, fewer false alarms)).
- Confusion matrix: TP=2085, FP=72, FN=1370, TN=56473 (FPR=0.13%, FNR=39.65%).
- KS is 0.6422 at threshold 0.12, indicating very strong separation.
- Top-10: precision=1.0000, recall=0.0029 shows the quality of the highest-risk shortlist.
- Top-500: precision=0.9620, recall=0.1392 shows how much coverage you get with a larger queue.



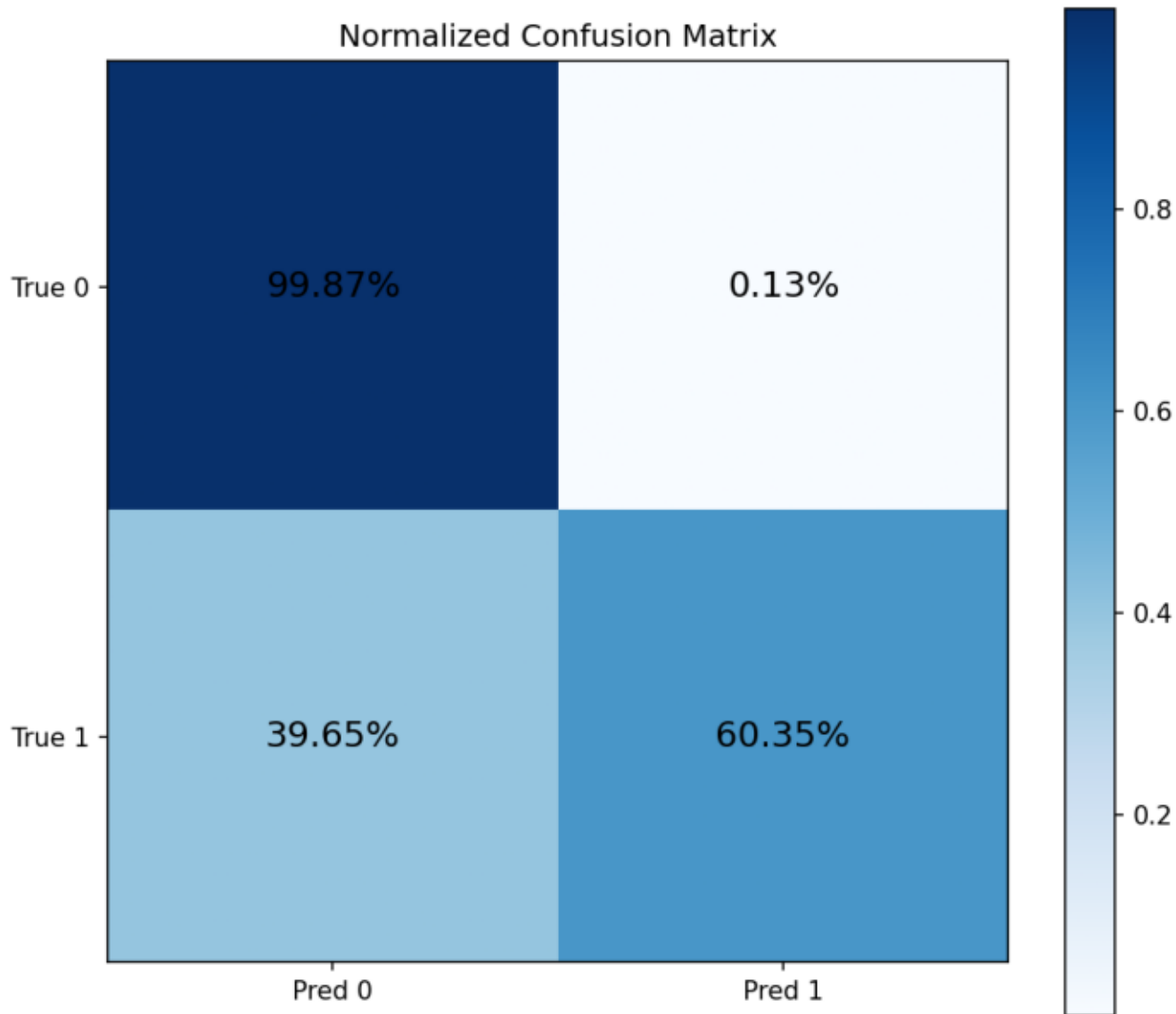
roc\_curve: AUC-ROC=0.8299; curve above diagonal indicates strong separation.



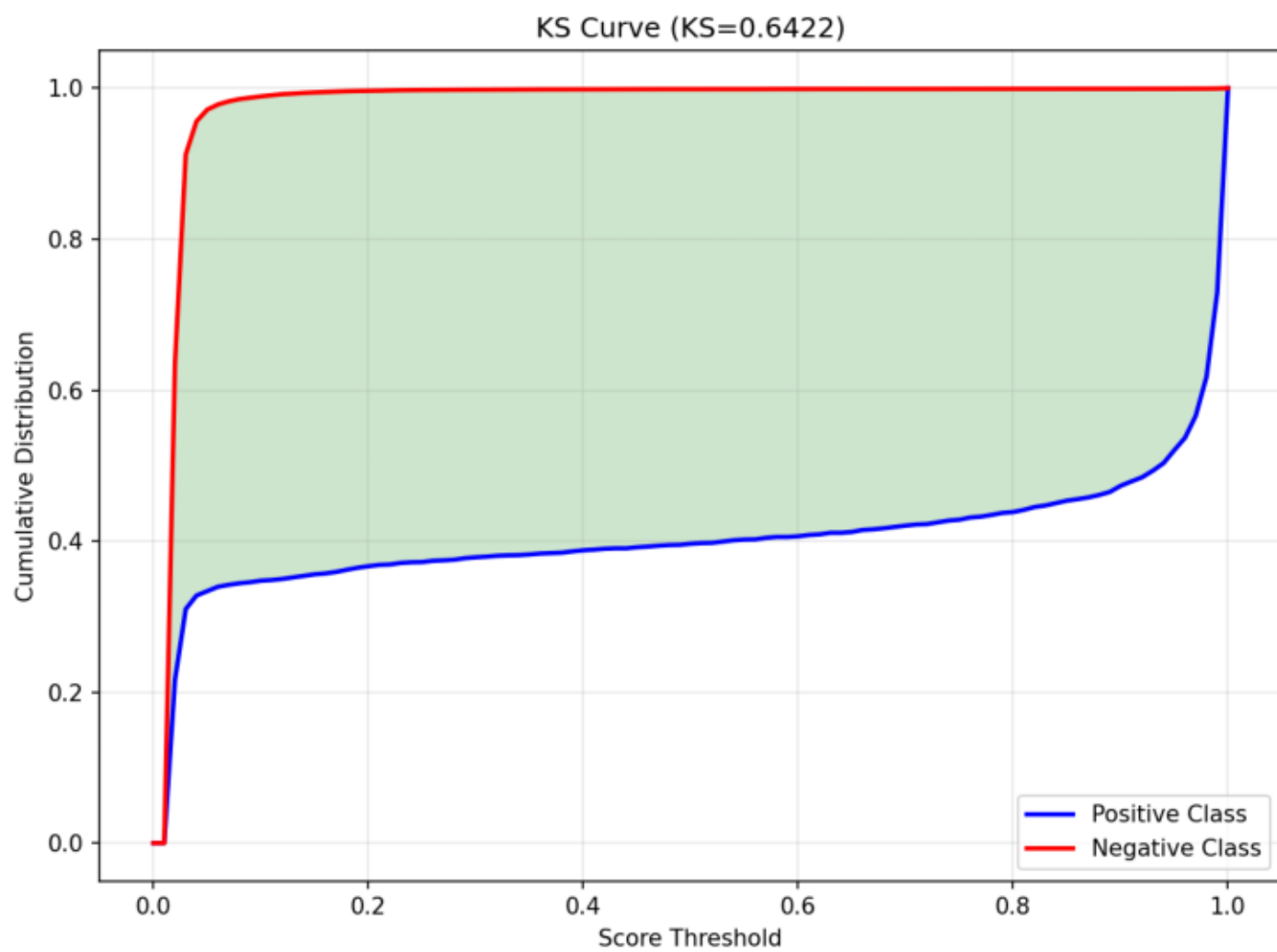
pr\_curve: AUC-PR=0.6833 vs baseline 5.76%; higher curve implies better precision at recall.



confusion\_matrix: At threshold 0.50: TP=2085, FP=72, FN=1370, TN=56473.

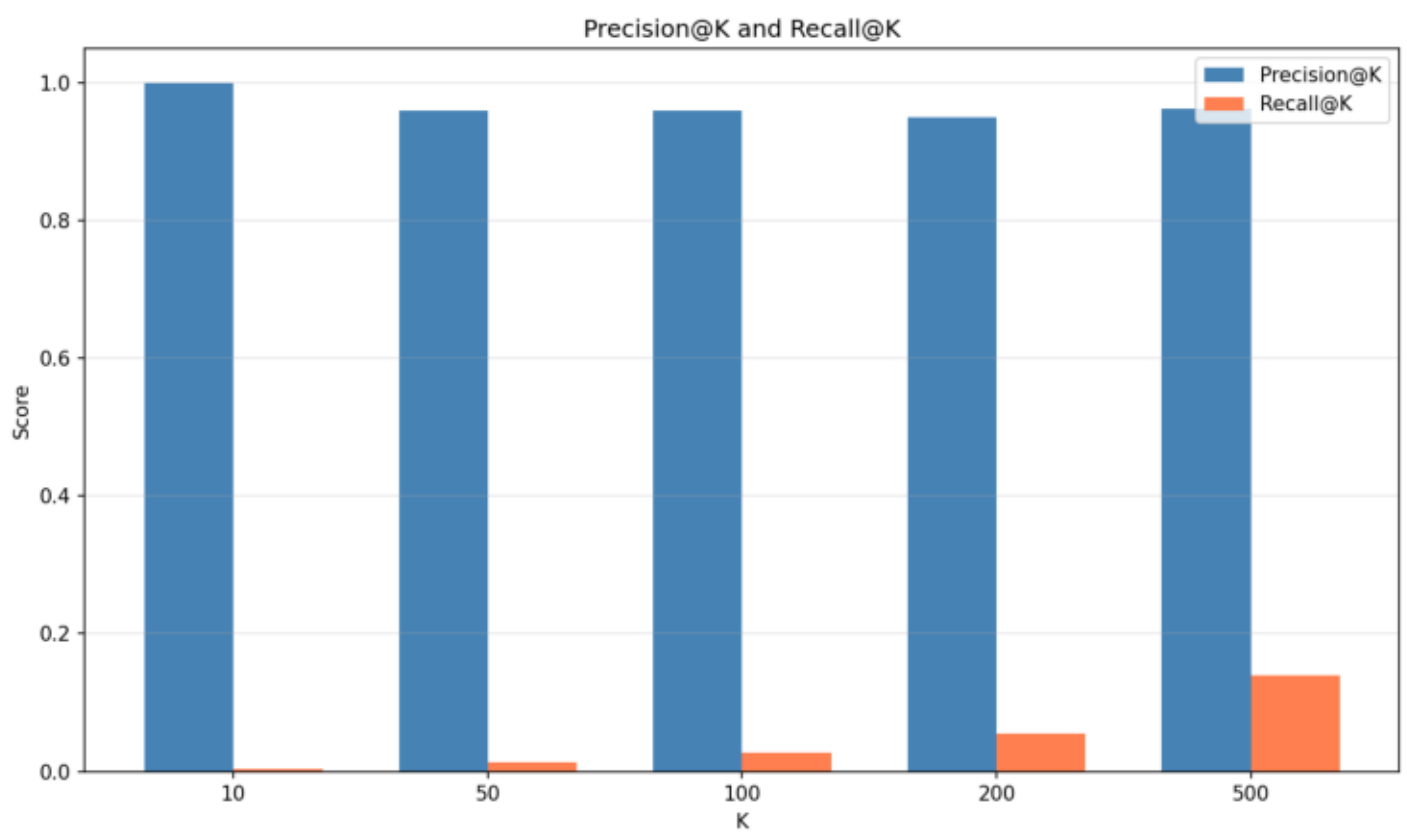


confusion\_matrix\_norm: Normalized rates: TPR=60.35%, TNR=99.87%.

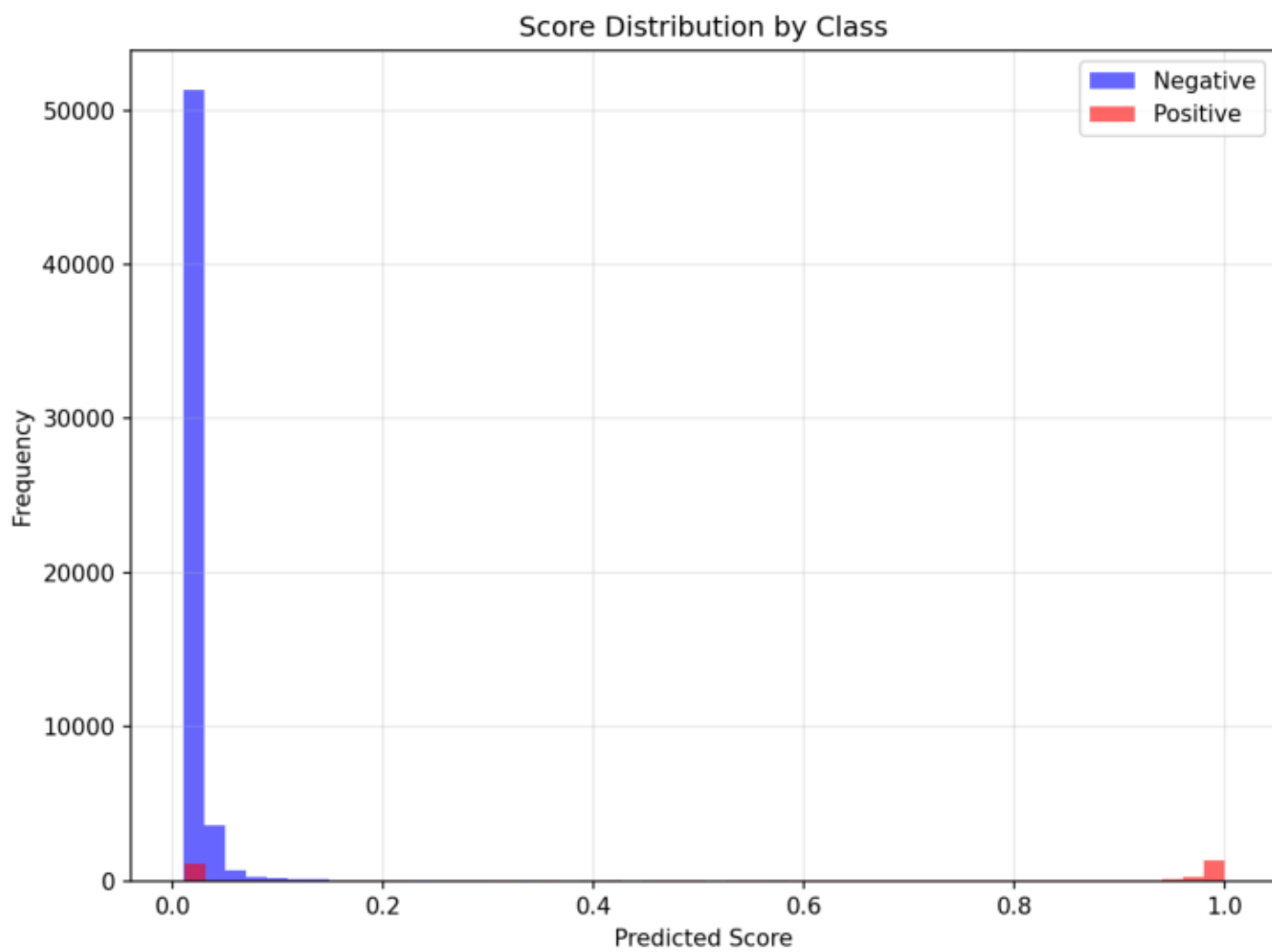


ks\_curve: Maximum separation KS=0.6422 at threshold 0.12.

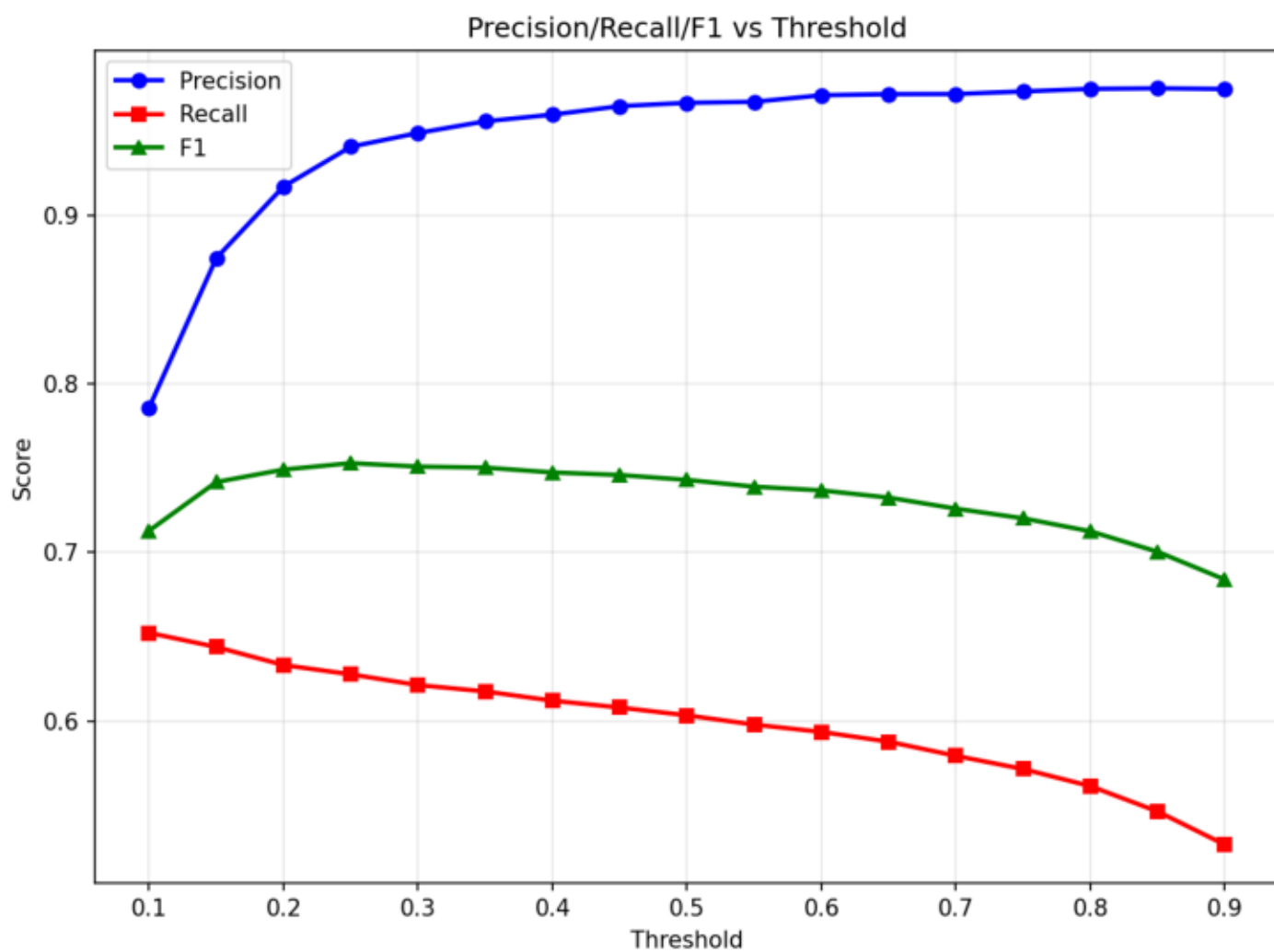




precision\_recall\_at\_k: Bars show precision/recall as you expand the review queue.



score\_distribution: Mean score: positive=0.597, negative=0.024 (larger gap is better).



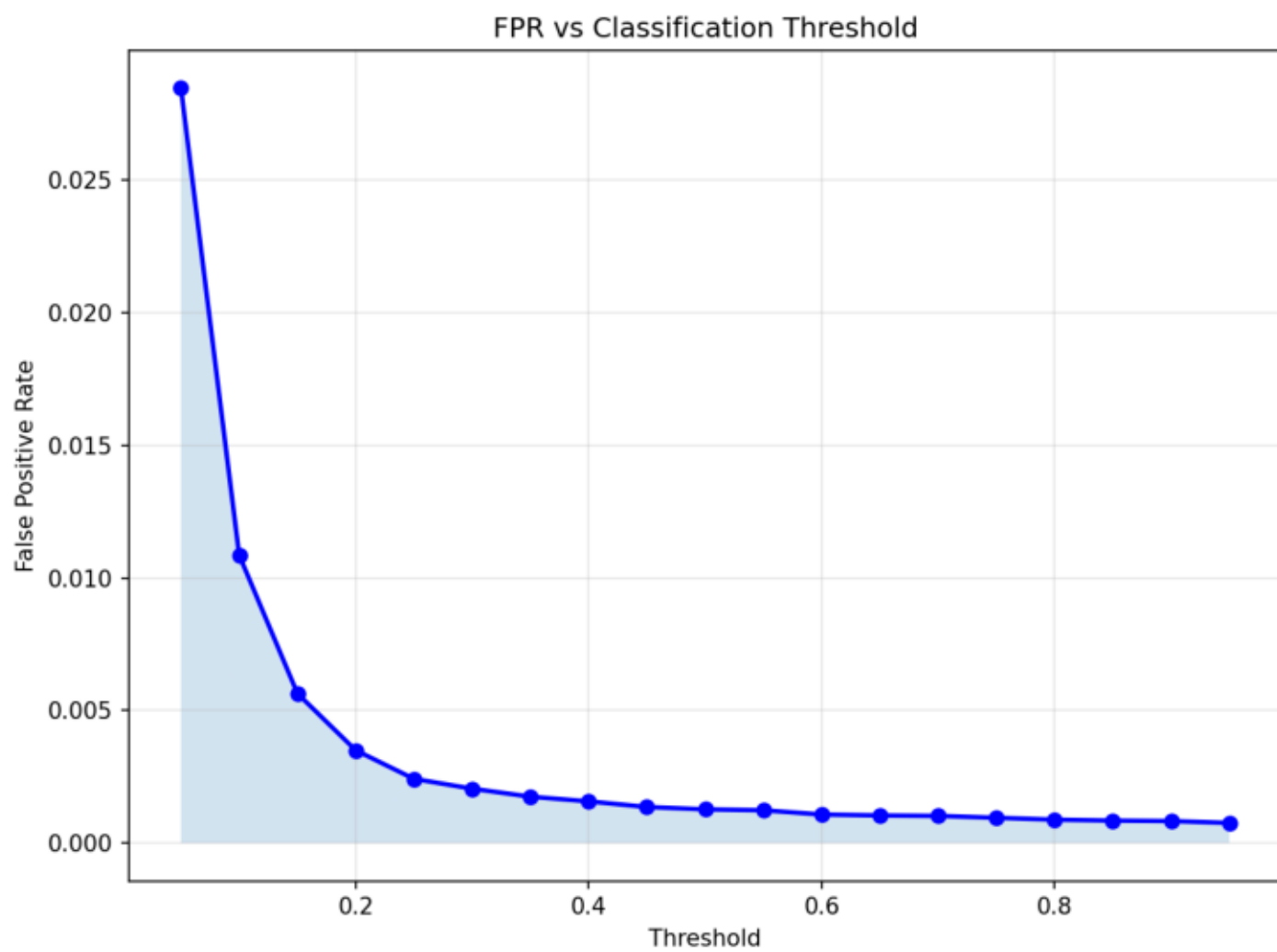
threshold\_analysis: Best F1 in tested grid is 0.7530 at threshold 0.25.

## 2) Model Efficiency

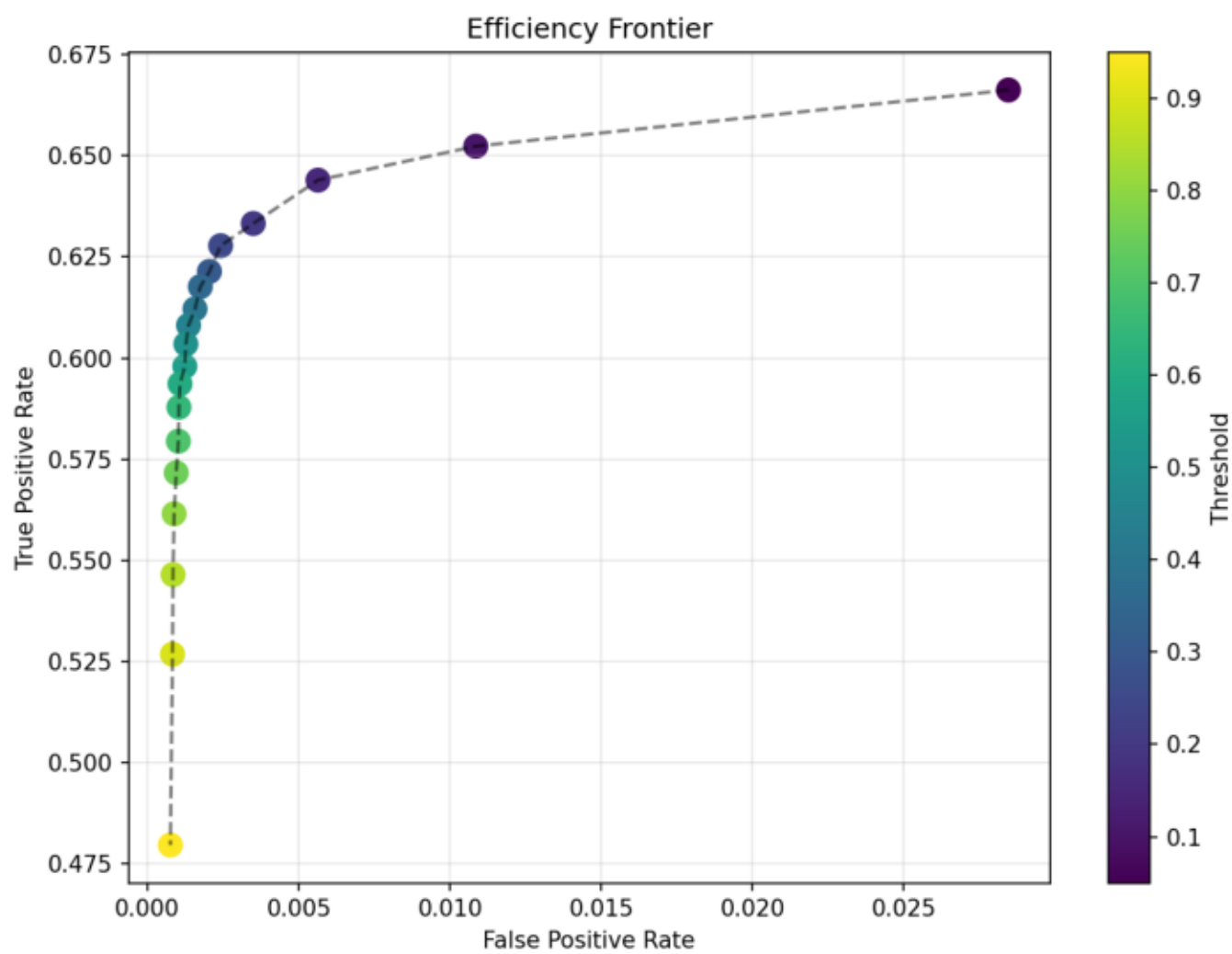
- fpr: 0.0013
- tn: 56473
- fp: 72
- threshold: 0.5000
- fpr\_at\_thresholds.t\_0.05: 0.0285
- fpr\_at\_thresholds.t\_0.10: 0.0109
- fpr\_at\_thresholds.t\_0.15: 0.0056
- fpr\_at\_thresholds.t\_0.20: 0.0035
- fpr\_at\_thresholds.t\_0.25: 0.0024
- fpr\_at\_thresholds.t\_0.30: 0.0021
- fpr\_at\_thresholds.t\_0.35: 0.0018
- fpr\_at\_thresholds.t\_0.40: 0.0016
- fpr\_at\_thresholds.t\_0.45: 0.0014
- fpr\_at\_thresholds.t\_0.50: 0.0013
- fpr\_at\_thresholds.t\_0.55: 0.0012
- fpr\_at\_thresholds.t\_0.60: 0.0011
- fpr\_at\_thresholds.t\_0.65: 0.0010
- fpr\_at\_thresholds.t\_0.70: 0.0010
- fpr\_at\_thresholds.t\_0.75: 0.0010
- fpr\_at\_thresholds.t\_0.80: 0.0009
- fpr\_at\_thresholds.t\_0.85: 0.0008
- fpr\_at\_thresholds.t\_0.90: 0.0008
- fpr\_at\_thresholds.t\_0.95: 0.0008

## 2) Model Efficiency - Explanations

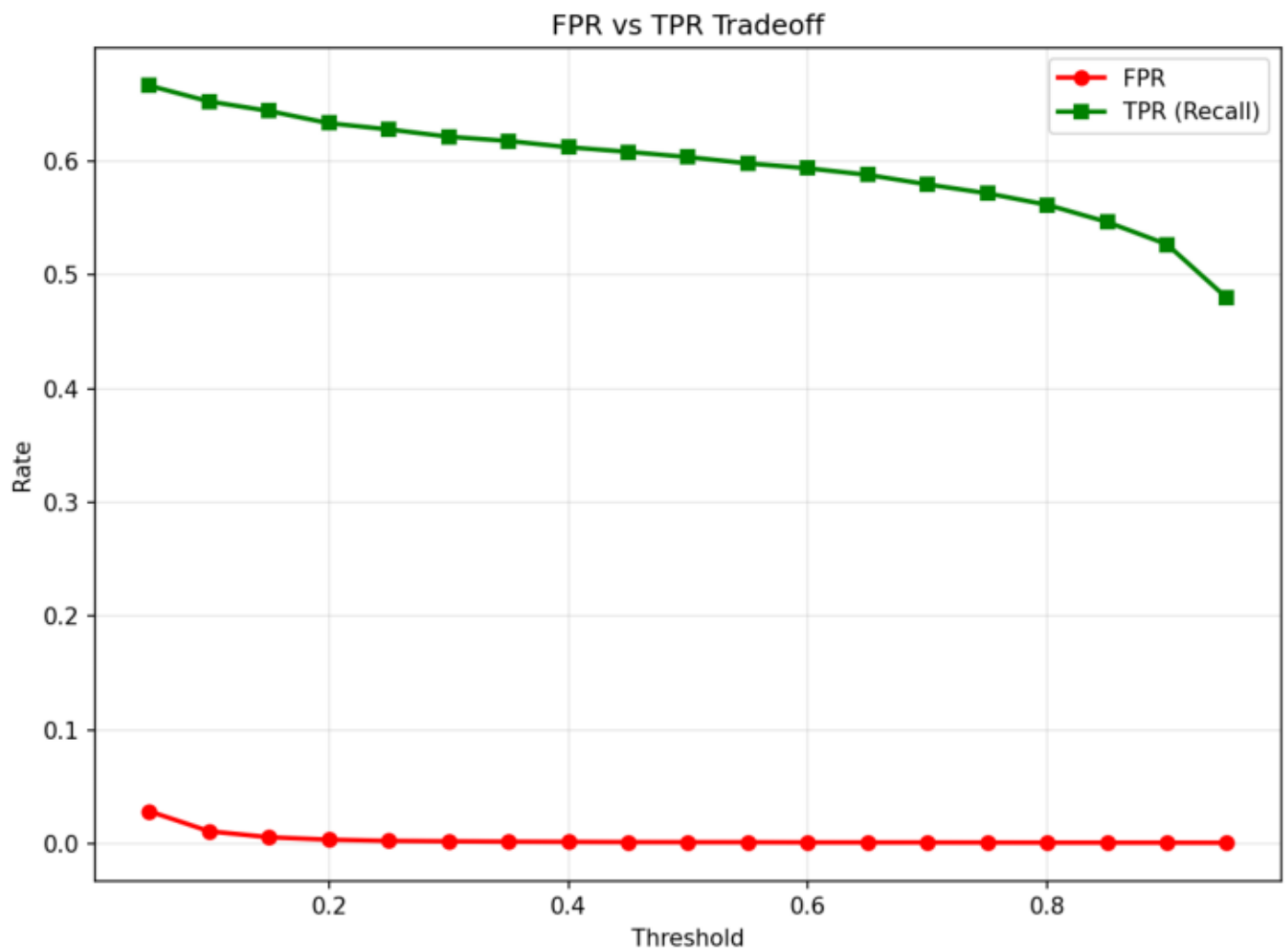
- At threshold 0.50,  $FPR=0.0013$  (low);  $FP=72$  out of 56545 negatives.
- At the same threshold,  $TPR=0.6035$  with  $TP=2085$  and  $FN=1370$ , showing the capture rate of positives.
- A threshold near 0.05 yields  $FPR \approx 0.0285$  with  $TPR \approx 0.6663$  if you want to target  $\sim 5\%$  false positives.



fpr\_vs\_threshold: FPR decreases as the threshold increases; use it to pick an operating point.



efficiency\_frontier: Each point shows the FPR/TPR tradeoff; move toward the top-left for better efficiency.



fpr\_tpr\_tradeoff: FPR and TPR curves highlight how recall drops as you reduce false positives.

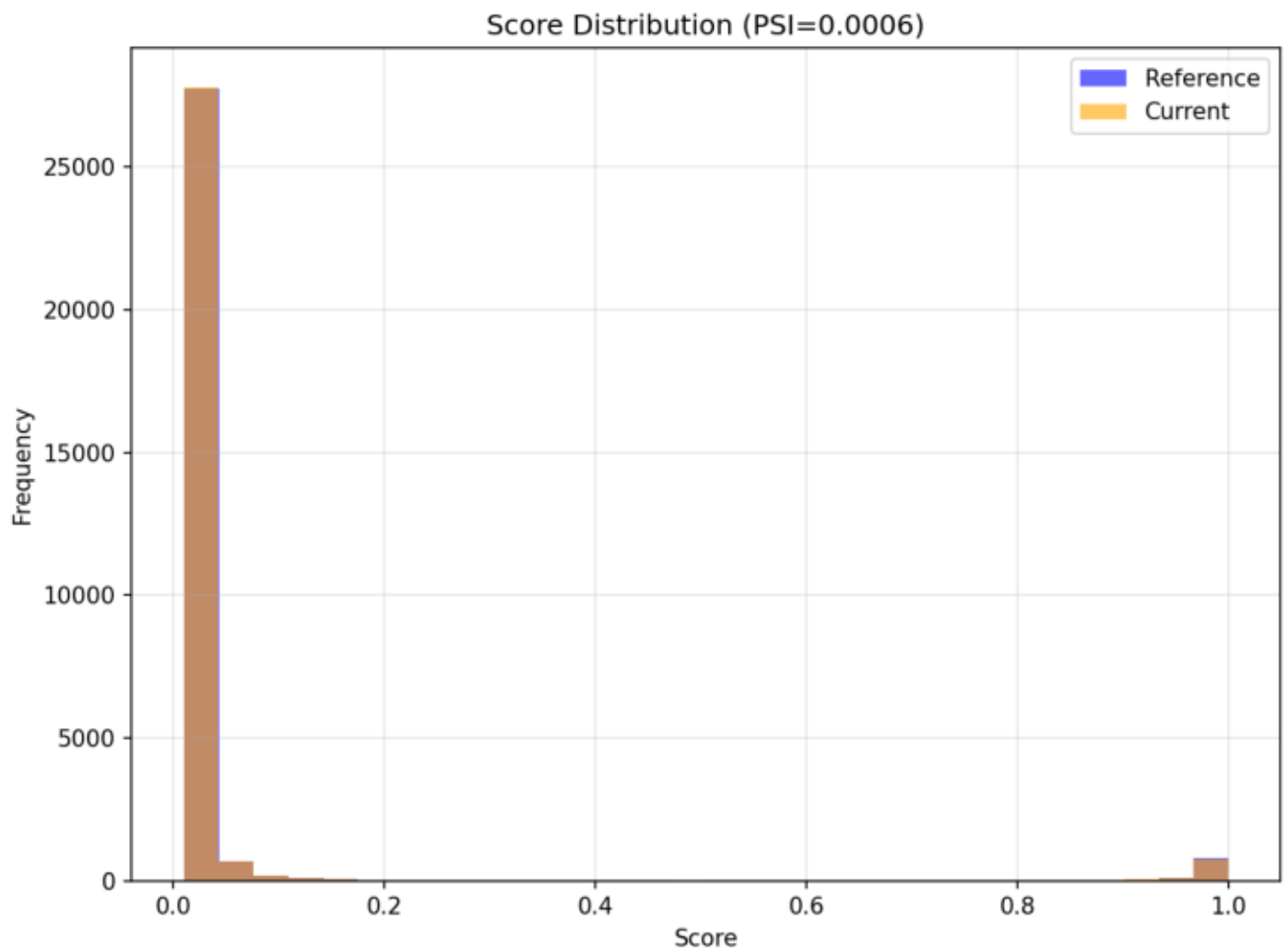


### 3) Model Stability

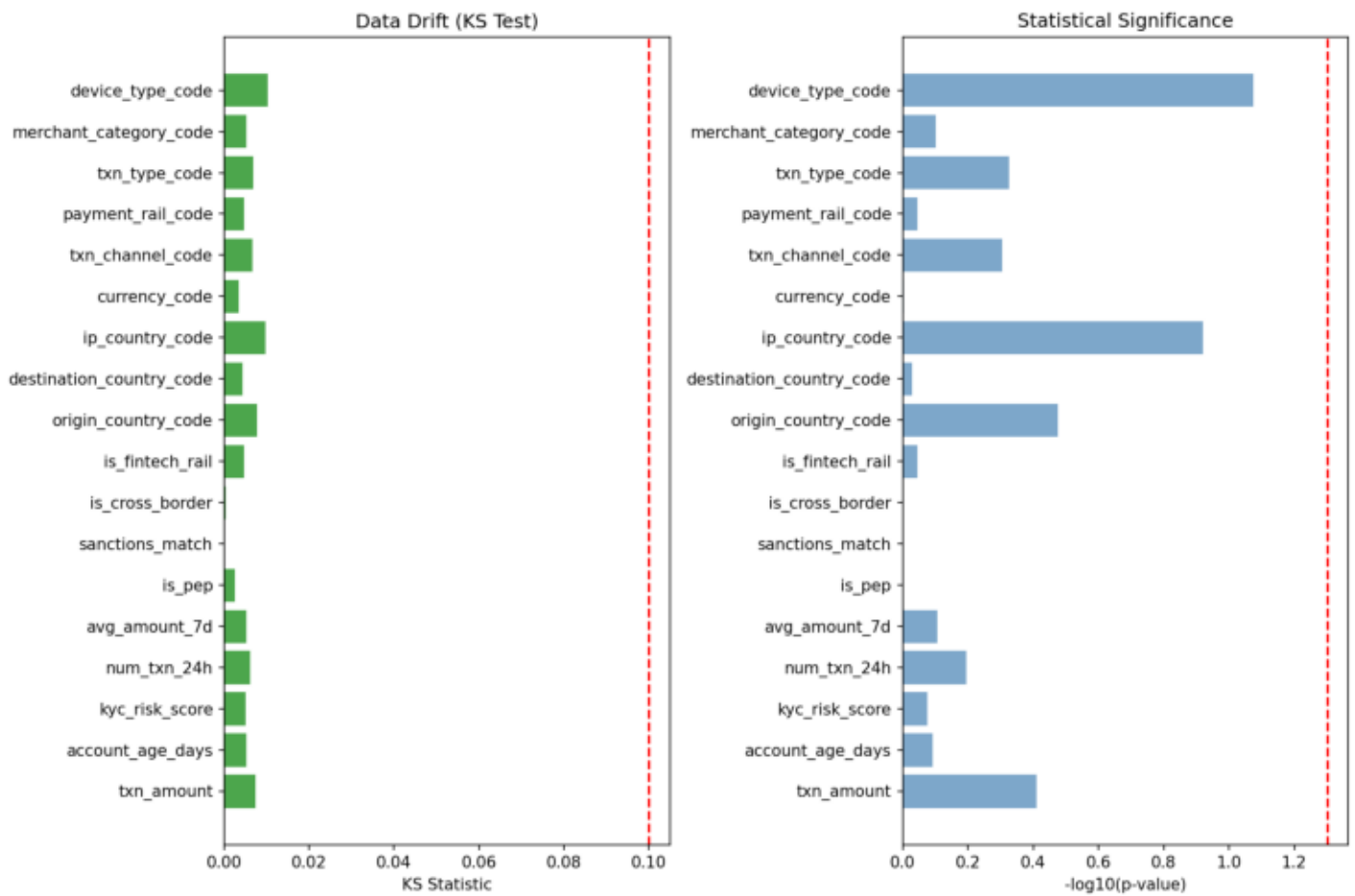
- psi: 0.0006
- cv\_auc\_roc\_mean: 0.8298
- cv\_auc\_roc\_std: 0.0107
- cv\_auc\_pr\_mean: 0.6837
- cv\_auc\_pr\_std: 0.0088
- bootstrap\_auc\_roc\_mean: 0.8357
- bootstrap\_auc\_roc\_ci\_lower: 0.8238
- bootstrap\_auc\_roc\_ci\_upper: 0.8467
- concept\_drift\_detected: False
- concept\_drift\_score: 0.0225

### 3) Model Stability - Explanations

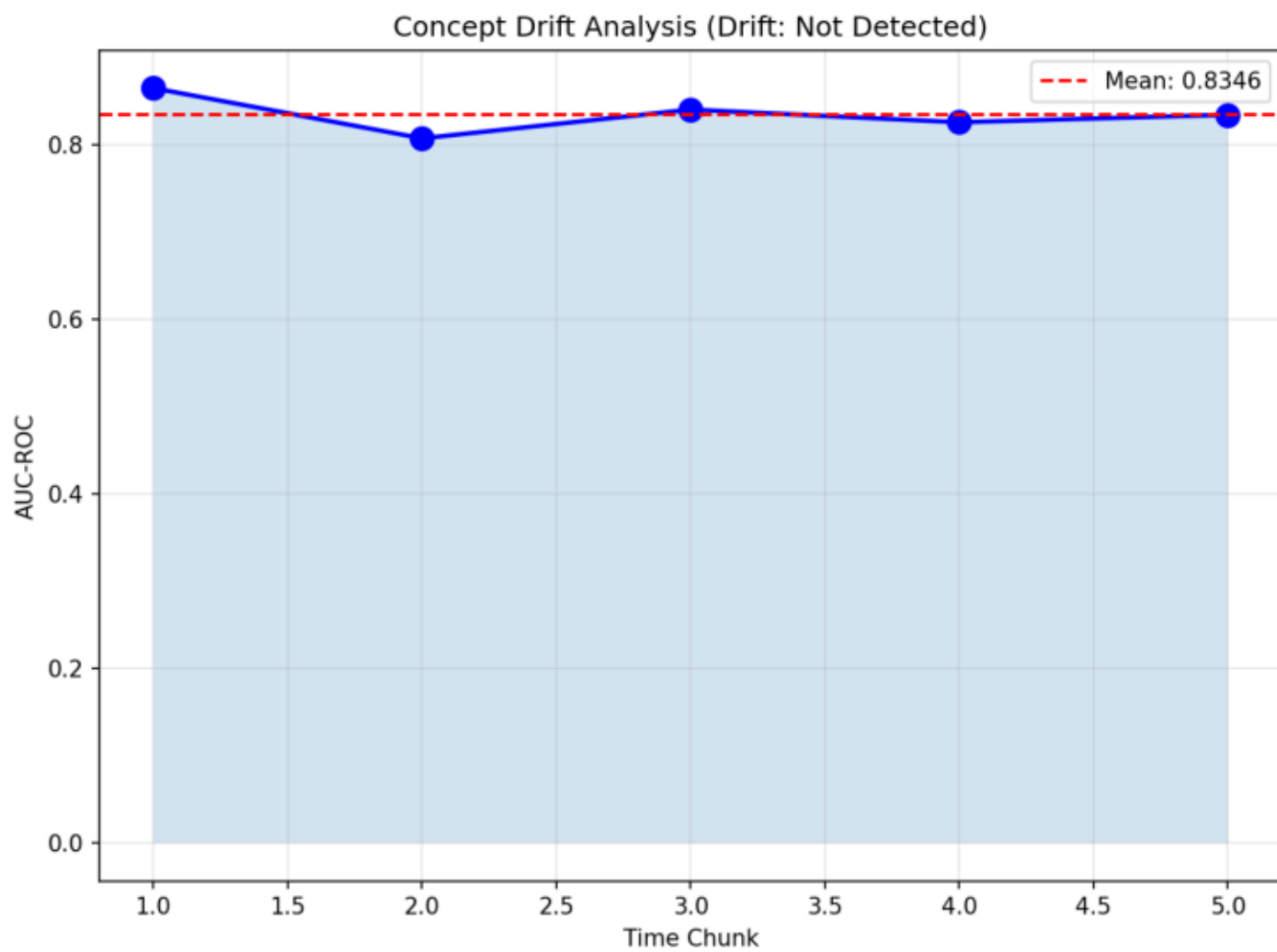
- PSI is 0.0006, indicating negligible score distribution shift between reference and current data.
- CV AUC-ROC is  $0.8298 \pm 0.0107$  (stable across folds).
- CV AUC-PR is  $0.6837 \pm 0.0088$  (stable across folds).
- Bootstrap AUC-ROC mean is 0.8357 with 95% CI [0.8238, 0.8467] (width 0.0230).
- No concept drift detected (score 0.0225), suggesting stable performance over time.
- Data drift flagged 0/18 numeric features; top drift signals: ['device\_type\_code', 'ip\_country\_code', 'origin\_country\_code'].



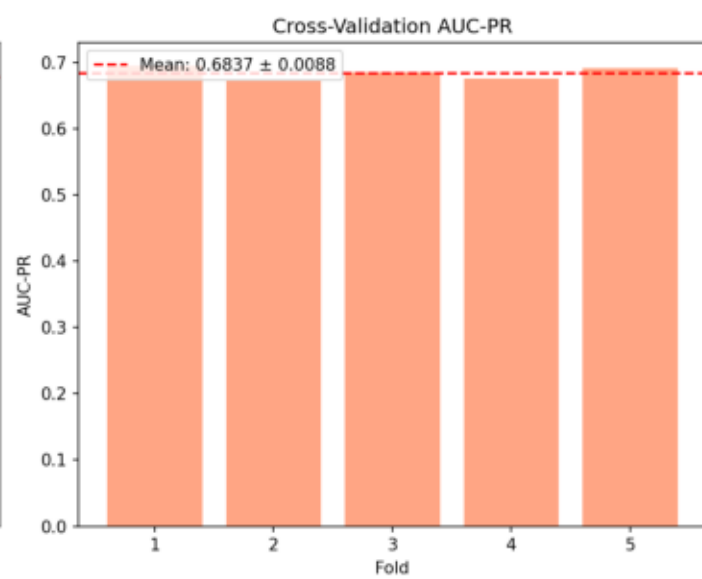
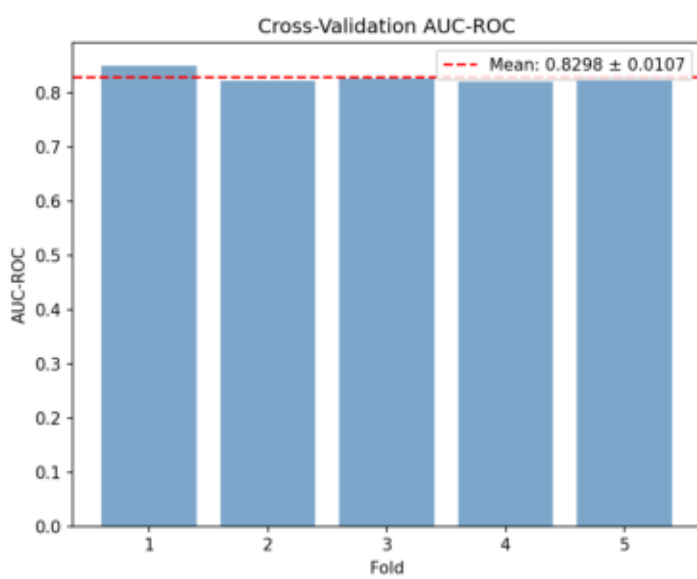
psi\_distribution: Score distributions overlap consistent with PSI=0.0006.



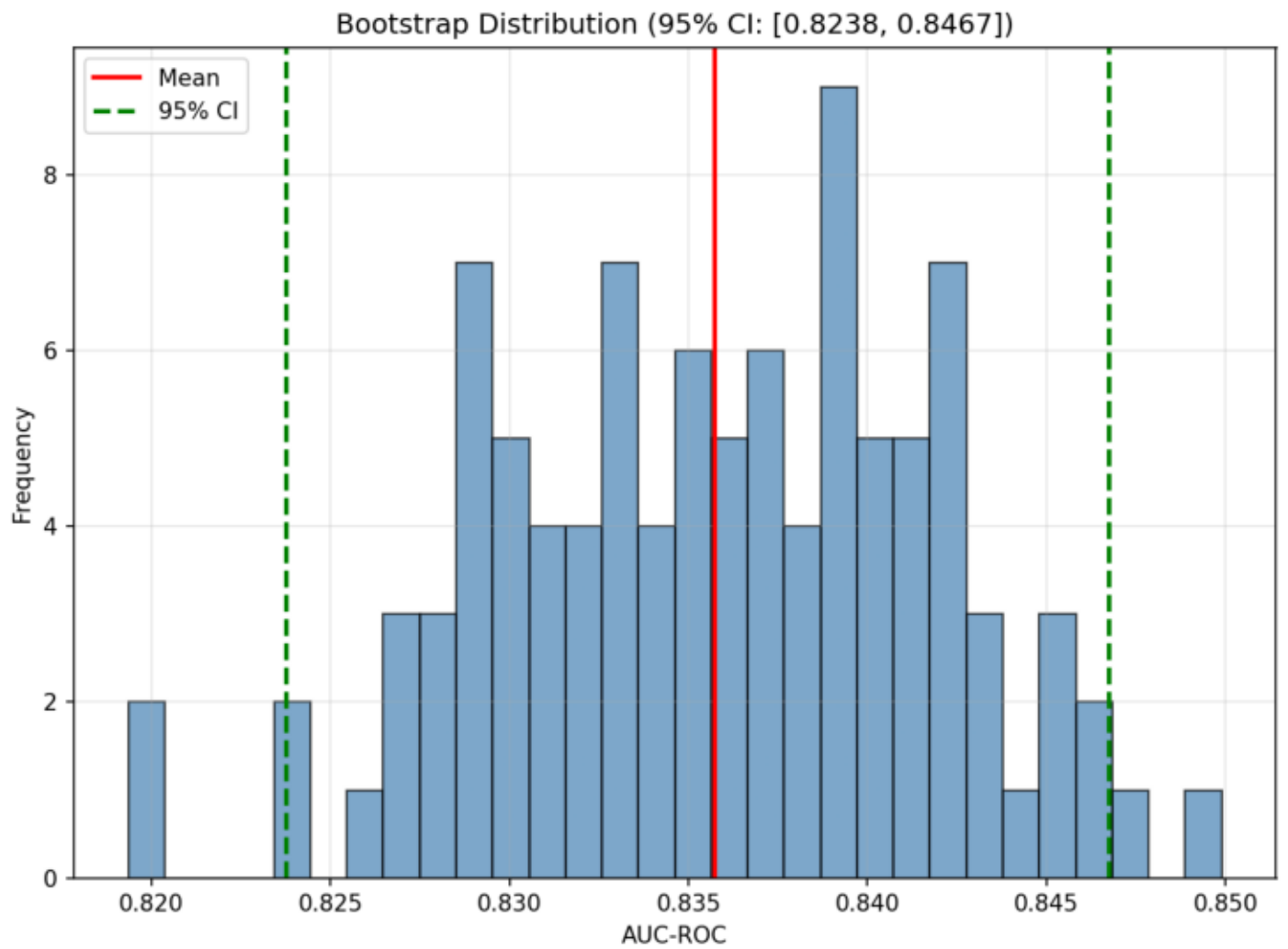
data\_drift\_heatmap: 0/18 features show drift; red bars highlight flagged features.



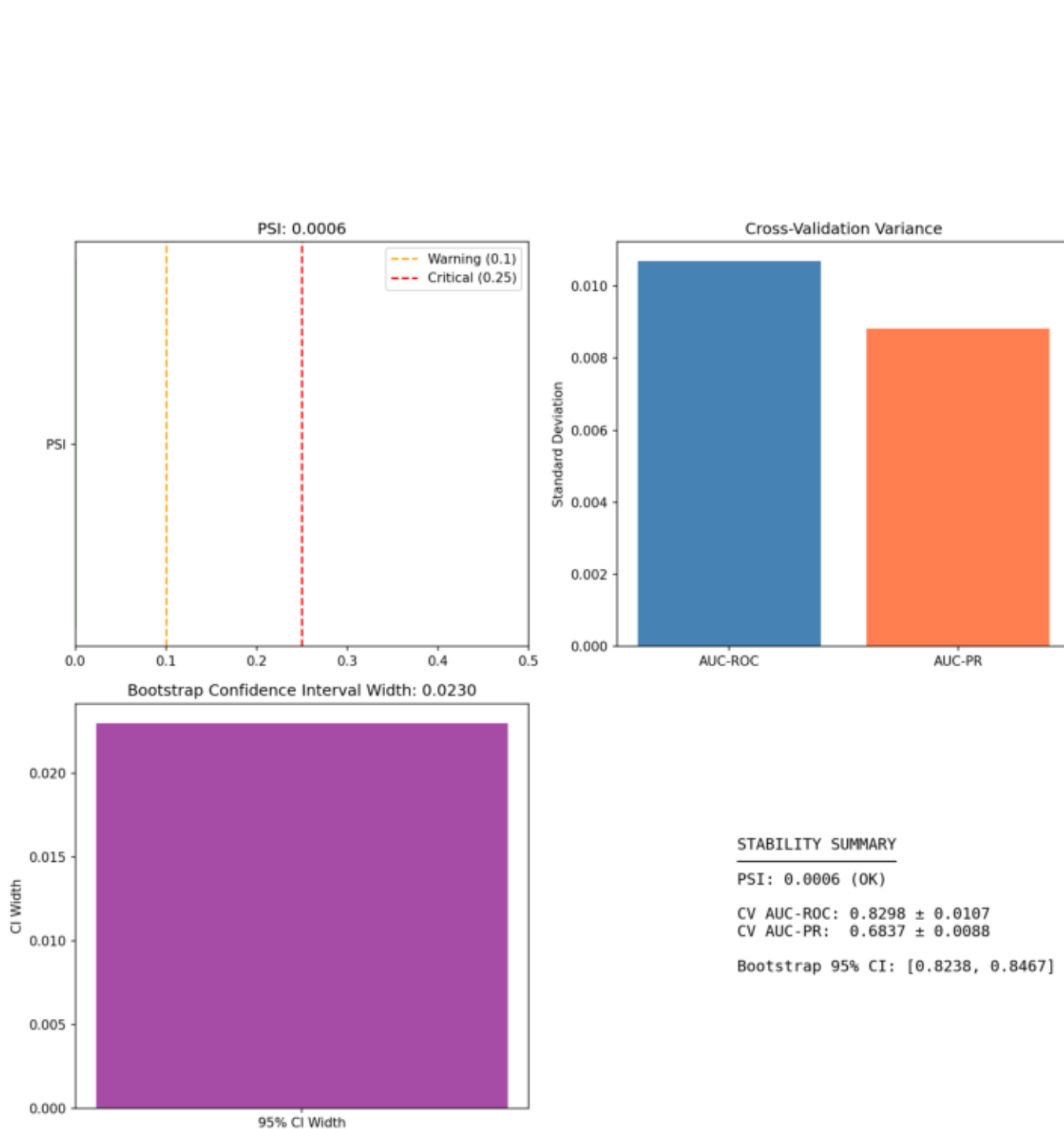
concept\_drift: AUC across time chunks shows whether performance is stable over time.



cv\_results: Fold scores cluster around ROC 0.8298 and PR 0.6837.



bootstrap\_distribution: Bootstrap CI [0.8238, 0.8467] reflects performance stability.



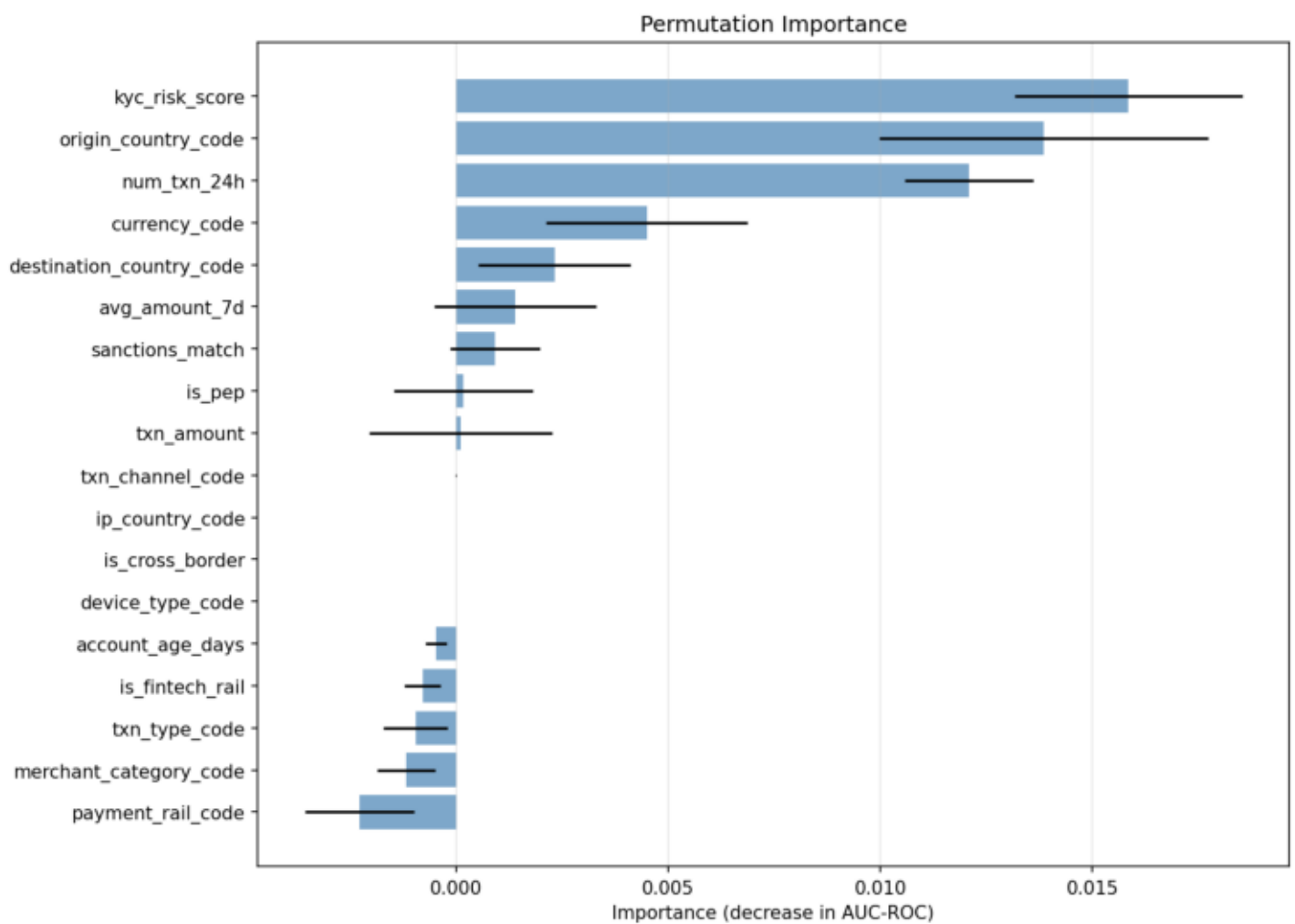


## 4) Model Interpretability

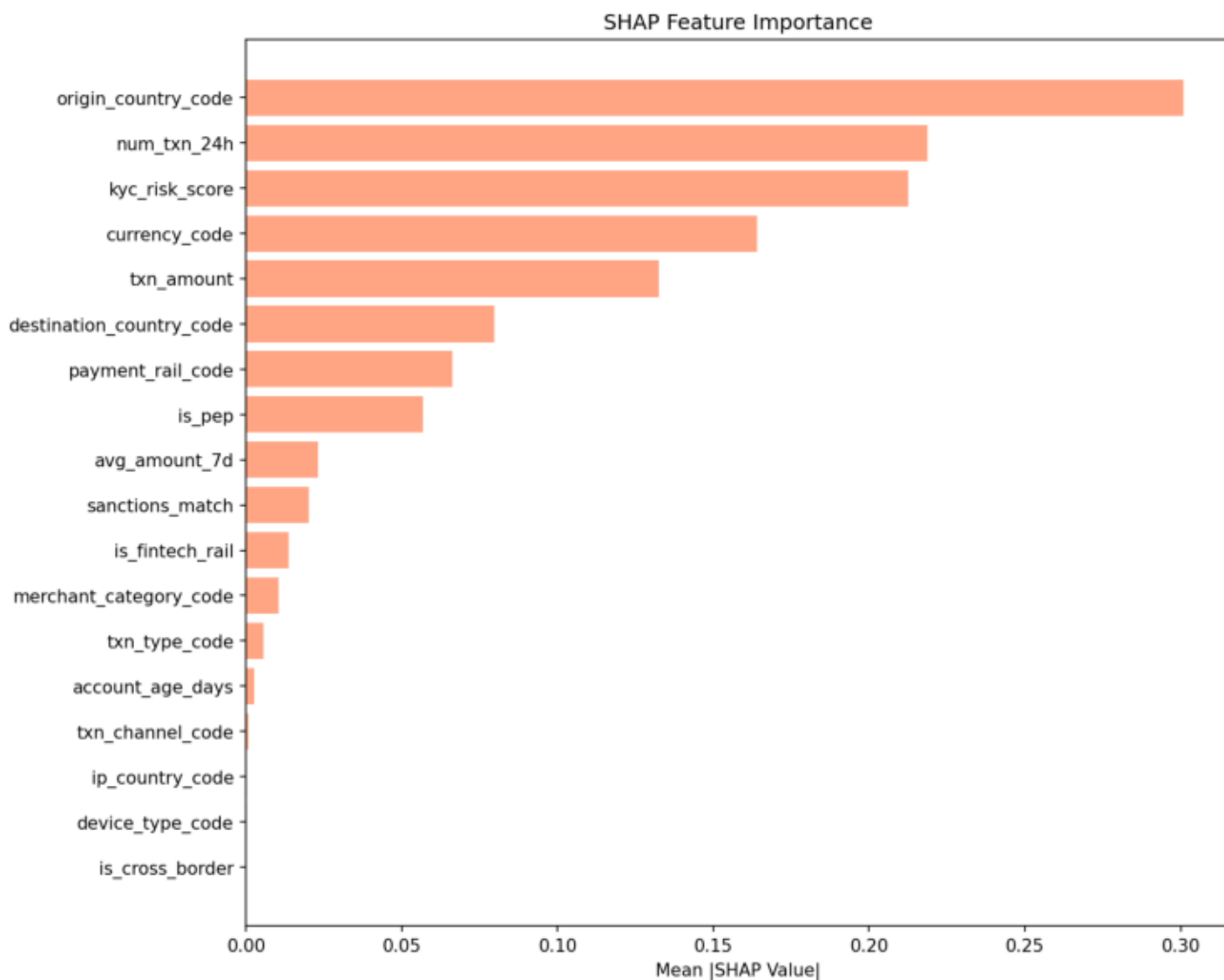
- model\_type: tree
- methods\_used: ['permutation', 'pdp', 'ice', 'shap', 'lime']
- perm\_top\_features: ['kyc\_risk\_score', 'origin\_country\_code', 'num\_txn\_24h', 'currency\_code', 'destination\_country\_code']
- shap\_top\_features: ['origin\_country\_code', 'num\_txn\_24h', 'kyc\_risk\_score', 'currency\_code', 'txn\_amount']...
- lime\_instances: 3
- pdp\_features: ['txn\_amount', 'account\_age\_days', 'avg\_amount\_7d', 'kyc\_risk\_score', 'num\_txn\_24h']...
- ice\_features: ['txn\_amount', 'account\_age\_days', 'avg\_amount\_7d', 'kyc\_risk\_score']

## 4) Model Interpretability - Explanations

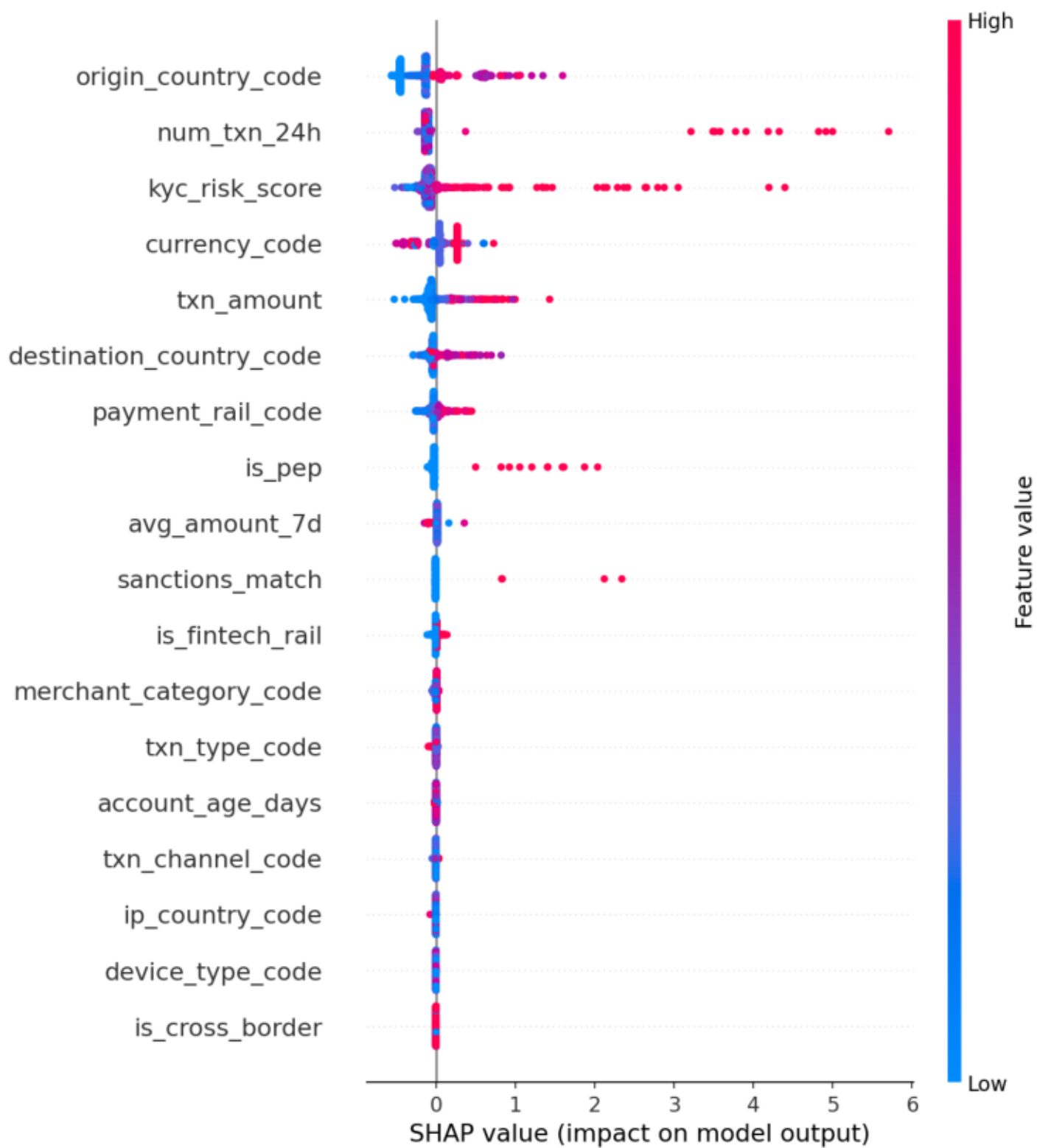
- Permutation importance highlights: ['kyc\_risk\_score', 'origin\_country\_code', 'num\_txn\_24h', 'currency\_code', 'destination\_country\_code'].
- SHAP top features: ['origin\_country\_code', 'num\_txn\_24h', 'kyc\_risk\_score', 'currency\_code', 'txn\_amount'].
- LIME generated local explanations for 3 instances; top contributors for instance\_6: ['num\_txn\_24h > 3.00', 'sanctions\_match <= 0.00', 'is\_pep <= 0.00', 'txn\_amount > 333.44', 'kyc\_risk\_score > 46.42'].
- PDP plots show average effects for: ['txn\_amount', 'account\_age\_days', 'avg\_amount\_7d', 'kyc\_risk\_score', 'num\_txn\_24h', 'origin\_country\_code'].
- ICE plots show individual-level effects for: ['txn\_amount', 'account\_age\_days', 'avg\_amount\_7d', 'kyc\_risk\_score'].



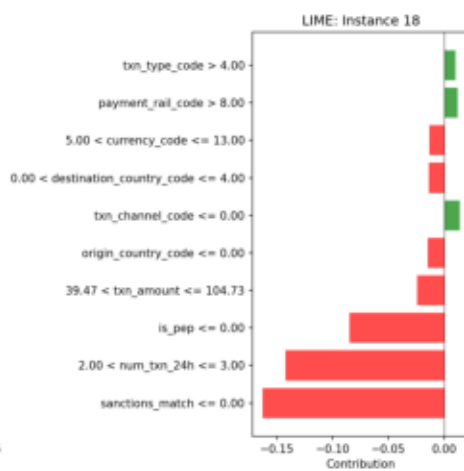
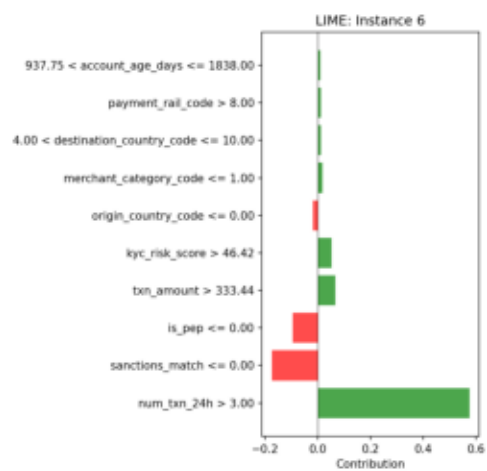
permutation\_importance: Bars show how much AUC drops when each feature is permuted.



shap\_bar: SHAP bar plot ranks global feature impact by mean absolute SHAP values.

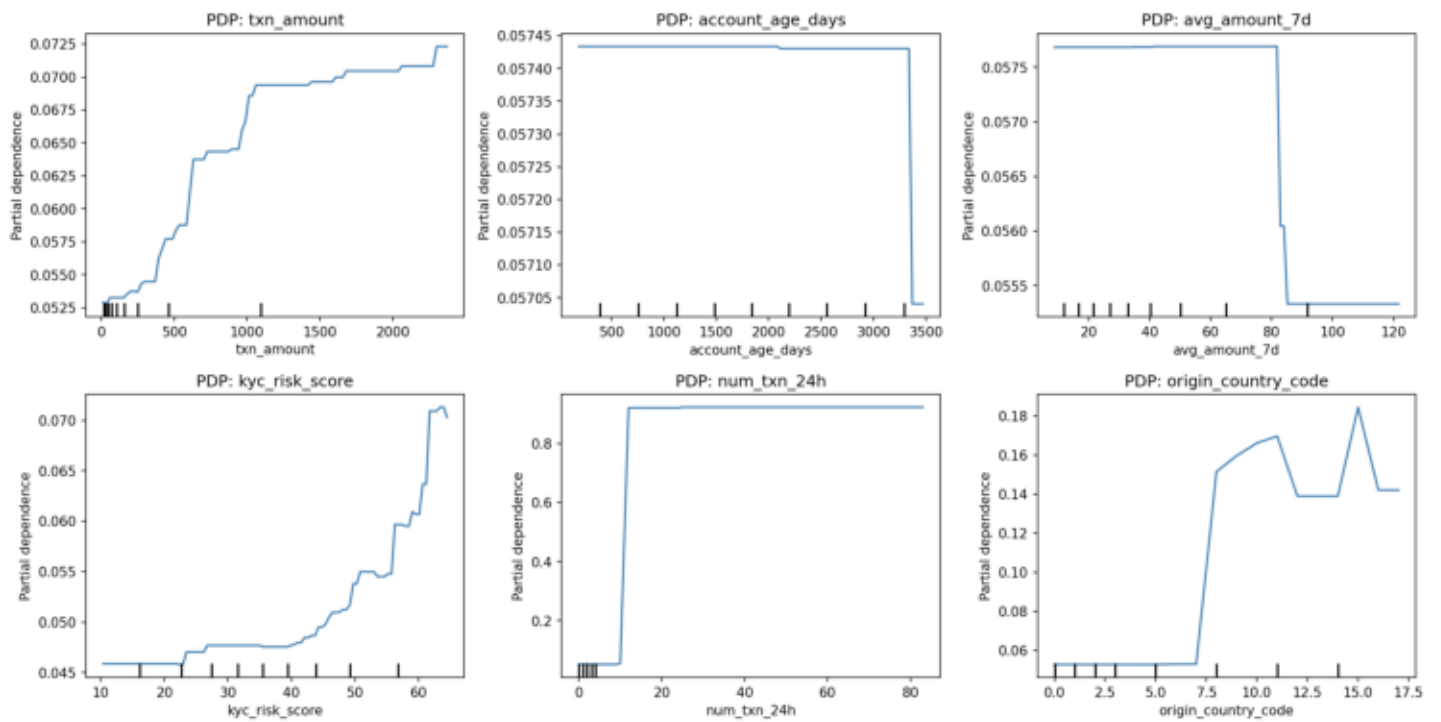


shap\_beeswarm: SHAP beeswarm shows both impact size and direction per feature.



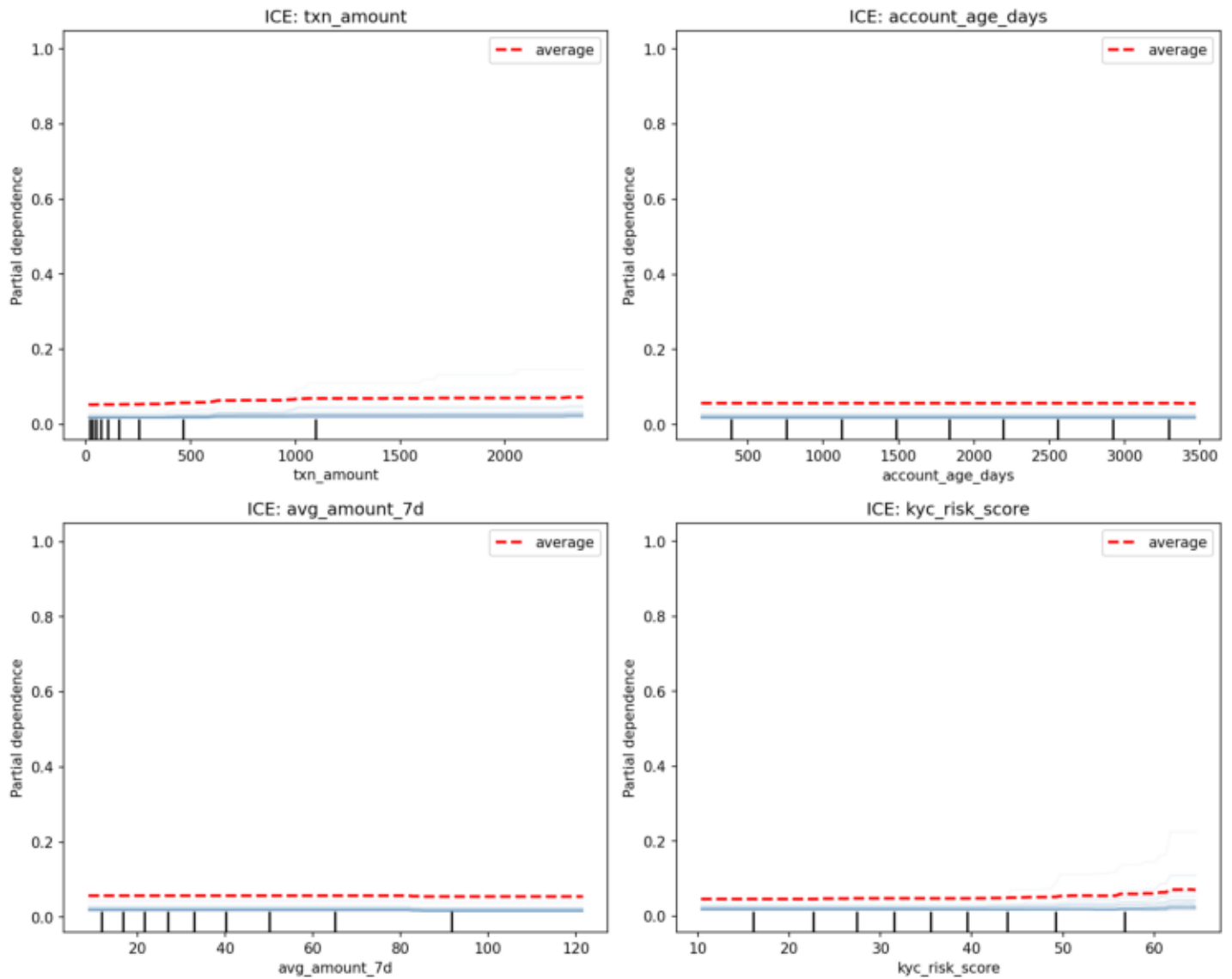
lime\_explanation: LIME bars show local feature contributions for selected instances.

Partial Dependence Plots



pdp: PDP curves show average model response as a feature varies.

### Individual Conditional Expectation (ICE) Plots



ice: ICE curves show per-instance responses as a feature varies.