

EXPLANATION LOTTERY: USER DECISION FRAMEWORK

INPUT: SHAP explanations from K models for instance x

STEP 1: Compute pairwise Spearman correlations between SHAP vectors

STEP 2: Calculate mean agreement (ρ) and standard deviation (σ)

STEP 3: Classify into decision category based on thresholds

Category	Condition	Recommended Action
✓ STRONG AGREEMENT (40.8% of cases)	$\rho \geq 0.7$	ACCEPT - Safe for automated use Explanation is reliable
~ WEAK AGREEMENT (23.2% of cases)	$0.5 \leq \rho < 0.7$	CAUTION - Document uncertainty Use but note limitations
△ MODERATE DISAGREE (17.4% of cases)	$0.3 \leq \rho < 0.5$	FLAG - Seek validation Additional review recommended
✗ STRONG DISAGREE (18.6% of cases)	$\rho < 0.3$	REJECT - Human review required Do NOT use single explanation

KEY FINDINGS:

- 36.0% of predictions need human review ($\rho < 0.5$)
- 40.8% are safe for fully automated decision-making ($\rho \geq 0.7$)
- Per Krishna et al. (2023): 86% of practitioners use ad-hoc heuristics
- Our framework reduces potential errors by 30.9%

REGULATORY IMPLICATIONS:

- EU AI Act Art. 13: 36.0% fail transparency requirements
 - Healthcare: 24.3% of clinical predictions unreliable
 - Finance: 24.3% of credit decisions may violate ECOA