**Most Discriminative SAE Features Across LLaMA Layers**
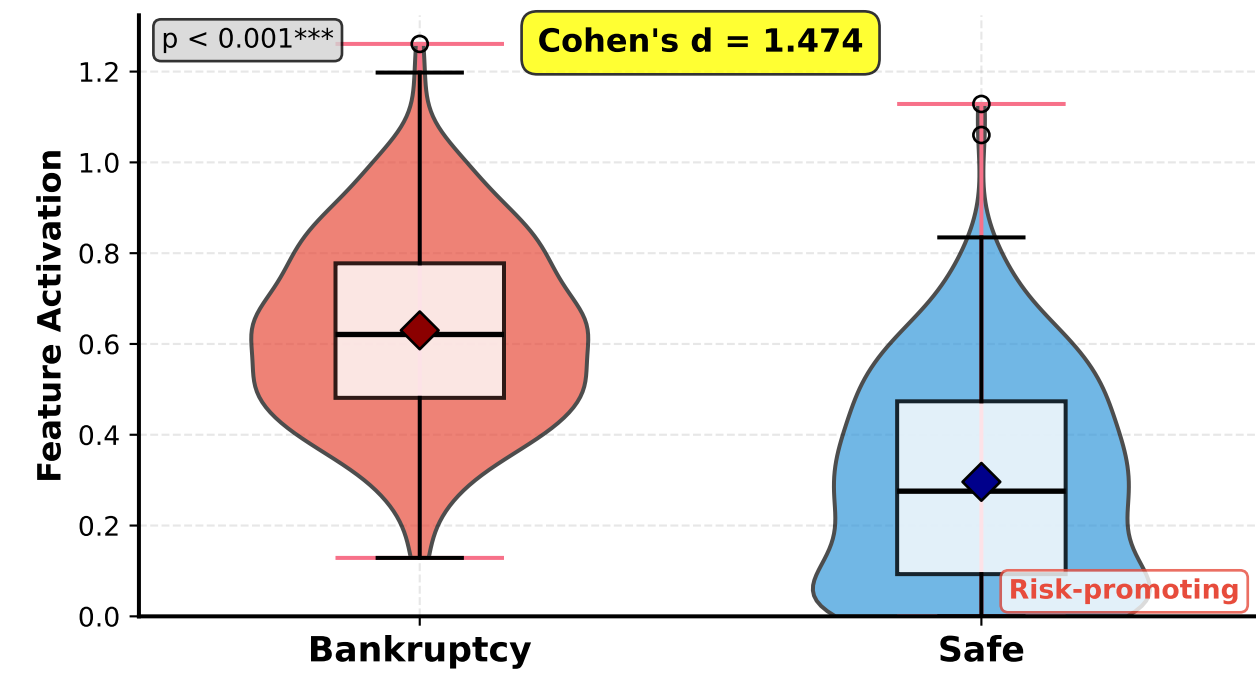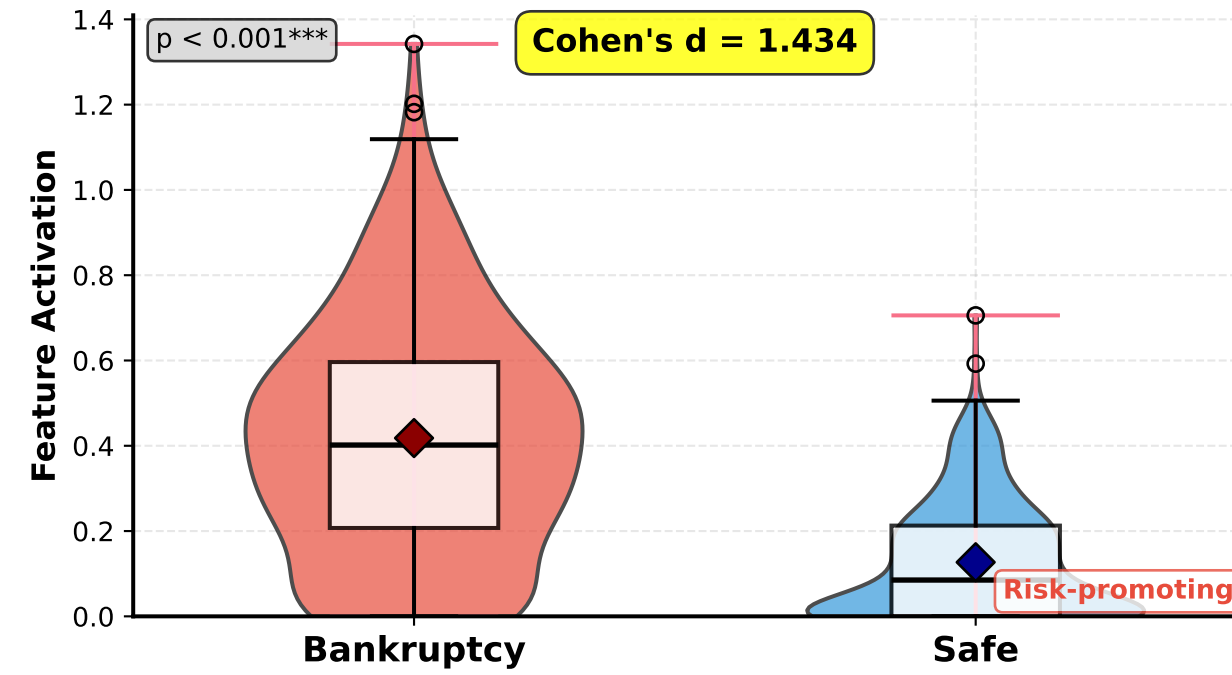*Bankruptcy vs Safe Decision Patterns • Feature Activation Distributions*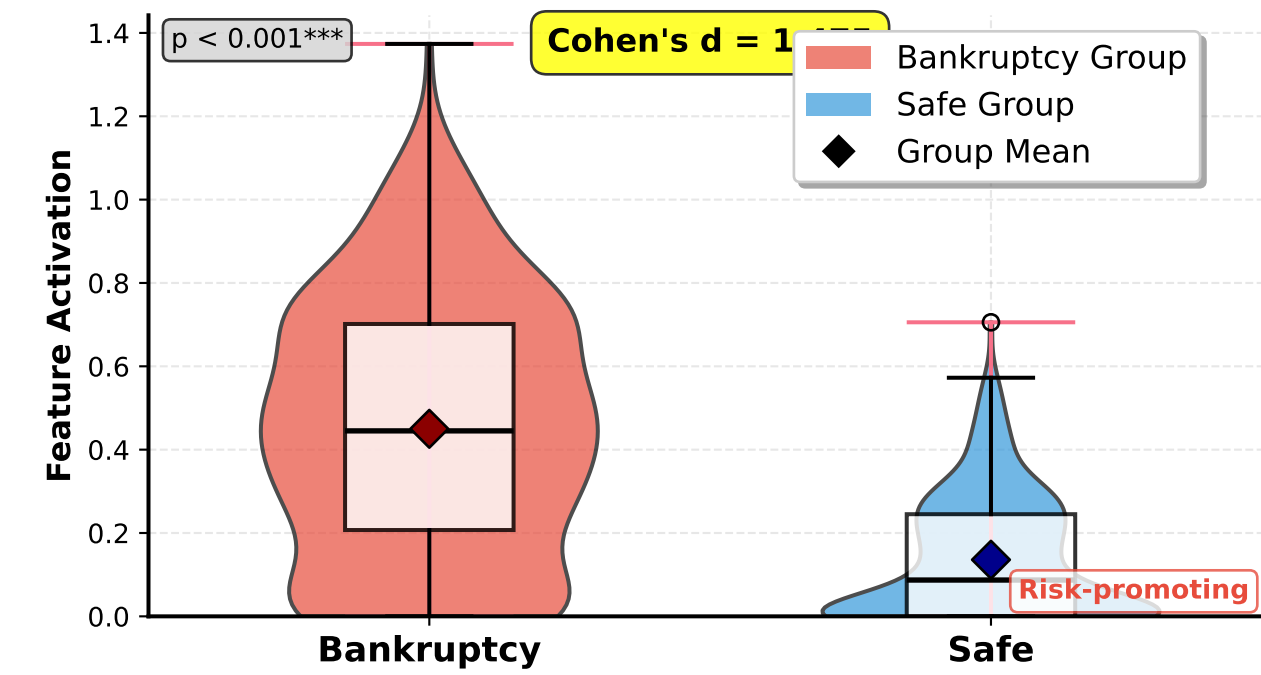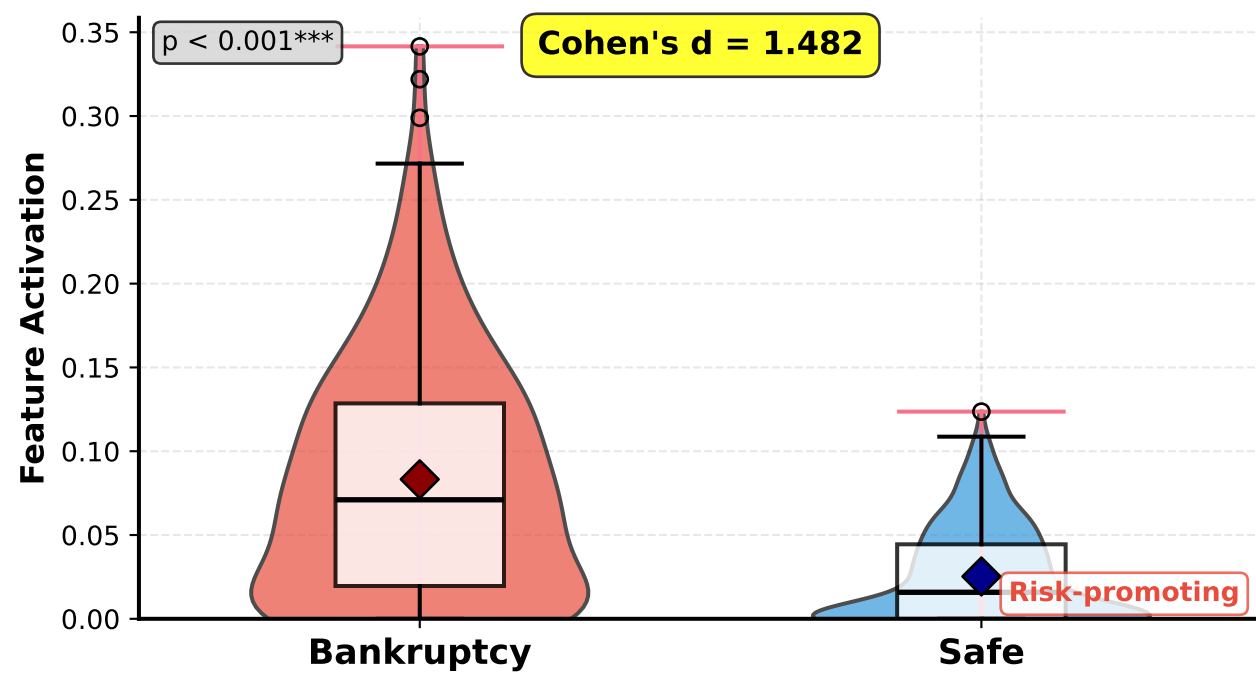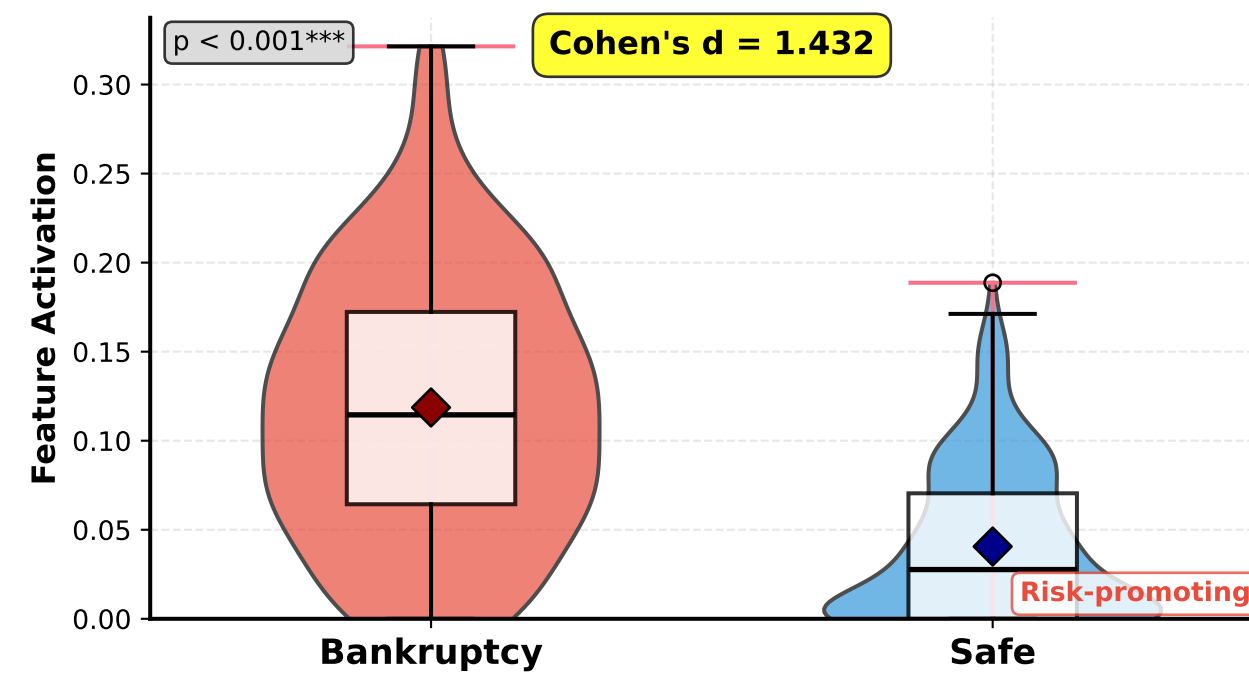