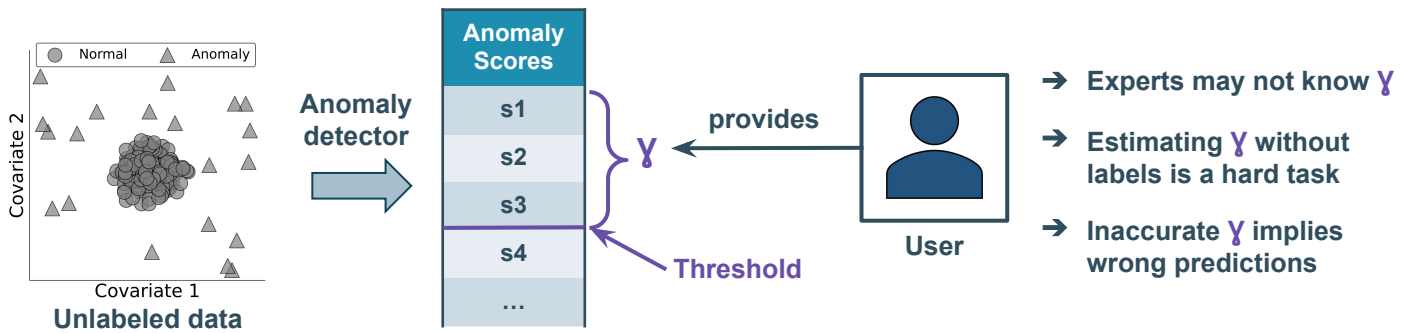
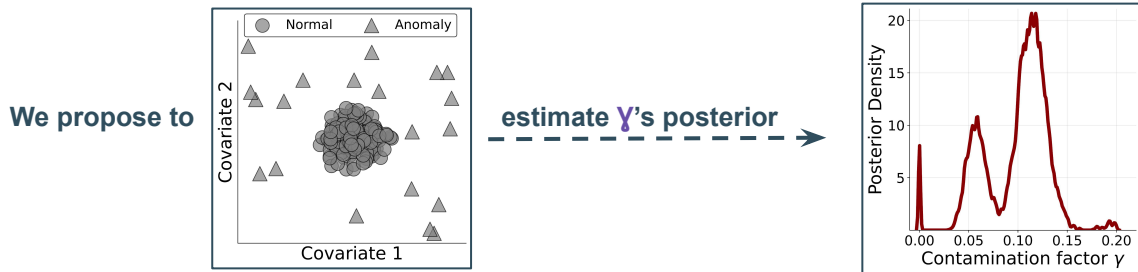


Estimating the Contamination Factor's Distribution in Unsupervised Anomaly Detection

Problem: thresholding the anomaly scores requires domain knowledge

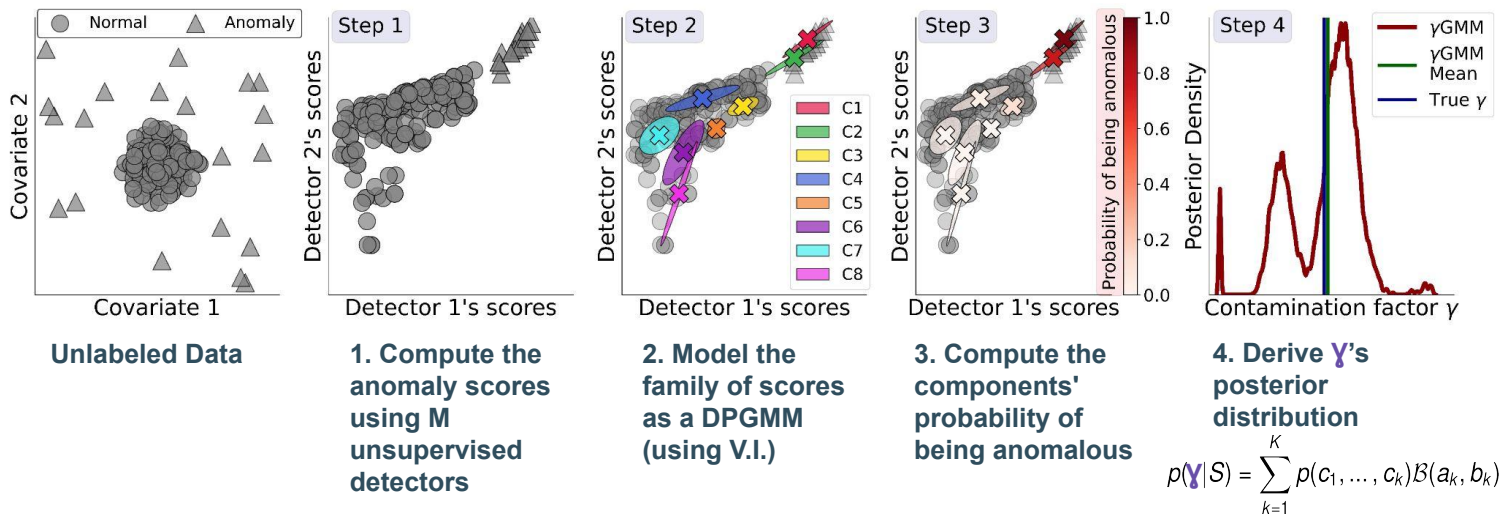


Task: estimate the contamination factor's posterior distribution using only unlabeled data



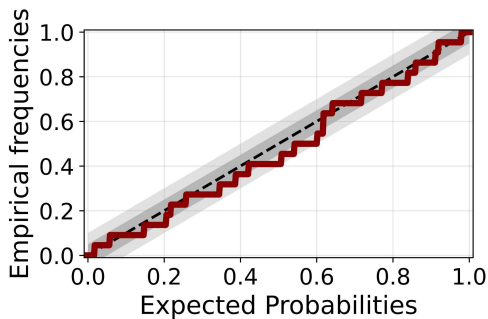
Insight: (a) Model the data in the anomaly score space, (b) identify the components flagged as anomalies by several detectors, and (c) estimate their mass as the contamination γ

We propose γ GMM, a 4-steps approach:

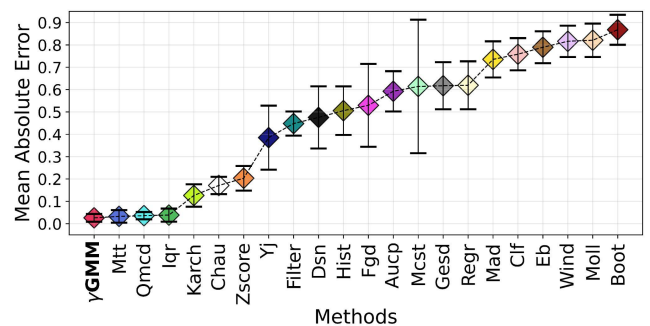


Experiments on 20 datasets show that γ GMM has

1. A well calibrated posterior



2. Low MAE when using the sample mean as point estimate



Lorenzo Perini, Paul Bürkner, Arto Klami

@LorenzoPerini95
<https://people.cs.kuleuven.be/~lorenzo.perini/>

