

Malware detection is a crucial aspect of cybersecurity, presenting several challenges, particularly in data stream scenarios that experience strong concept drift and label delay. The concept drift is characterized by the presence of highly influential yet rapidly changing features, such as specific filenames or mutexes, alongside stable features, such as connection types or monetization methods, which remain relatively consistent over time. In this thesis, we formalize this scenario and further exploit the hypothesis that the adaptive removal of severely drifting subsets of features may have a great impact on procedure performance. We indeed demonstrate that current methods exhibit shortcomings connected with these features, especially during short periods following the arrival of a new concept. To validate the hypothesis of performance improvement through adaptive feature elimination, we propose two solutions: one based on Hellinger distance concept drift detection and the other on an incremental Gaussian Mixture Model algorithm. We evaluate both approaches using real-life data and our synthetic dataset, showing significant improvements on the synthetic dataset and promising results on real-life data. Additionally, we provide a comprehensive explanation of the techniques employed in the thesis.