

Precision Evaluation of an Automated Issue Signal Mining Pipeline

Denario

Anthropic, Gemini & OpenAI servers. Planet Earth.

Automated detection of operational issues and user experience friction is crucial for maintaining system health and user satisfaction. This study evaluates the precision of an automated signal mining pipeline designed to identify such problems. We rigorously assessed 20 sampled problem-class signals from a dataset of 121 system rollups using human evaluation by an expert (Claude) to determine their validity. The pipeline achieved an overall precision of 75.0%, indicating that three out of four detected signals correspond to genuine issues. Class-specific analysis revealed a robust 100% precision for defect detection, while user experience friction signals showed 69% precision, suggesting a higher false positive rate in this category. Most confirmed issues were of medium severity (80.0%), and detailed false positive analysis attributed these errors primarily to insufficient contextual information or ambiguous signal interpretations. These findings offer directional insights into the pipeline’s efficacy, underscoring its strong capability in defect identification and outlining clear areas for enhancing UX friction signal accuracy.

I. INTRODUCTION

Modern software systems are fundamental to nearly every aspect of daily life and business operations, characterized by their immense scale, intricate interdependencies, and continuous evolution. Ensuring their seamless operation and delivering an optimal user experience (UX) are paramount for maintaining system health, user satisfaction, and business continuity. However, the sheer volume and velocity of operational data—ranging from system logs and performance metrics to user interactions and feedback—render manual monitoring and problem detection increasingly intractable. The timely identification of operational issues, such as software defects, infrastructure failures, or performance bottlenecks, as well as subtle forms of user experience friction, is critical for proactive intervention and minimizing downtime or user frustration.

The principal challenge lies in effectively sifting through this deluge of information to distinguish genuine problems from benign events or mere noise. Traditional rule-based monitoring systems often prove insufficient for the dynamic nature of complex systems, frequently leading to a high rate of false positives or, conversely, missing critical issues entirely. This difficulty is compounded by the often ambiguous nature of system signals, where a single data point might possess multiple interpretations, and the necessary contextual information for accurate diagnosis is frequently distributed across disparate data sources. Furthermore, the definition of what constitutes a “problem” can be inherently subjective, particularly within the realm of user experience, making automated detection a complex endeavor. Consequently, the development of robust, automated solutions for issue detection has become a central focus in site reliability engineering and product management.

To address this pressing need, automated signal mining pipelines have emerged as a promising approach. These pipelines are specifically designed to leverage advanced analytical techniques, often incorporating machine learning and anomaly detection algorithms, to automatically identify patterns indicative of underlying operational problems and UX friction. By processing vast and diverse datasets, these systems aim to provide early warnings and actionable insights, thereby reducing the Mean Time To Resolution (MTTR) and enhancing overall system resilience. However, the true utility of such a pipeline hinges critically on its ability to accurately identify genuine issues without inundating human operators with irrelevant or misleading alerts.

A high rate of false positives can erode trust in an automated system, lead to alert fatigue among operators, and ultimately undermine the very goal of automation, turning a beneficial tool into an operational burden. Therefore, this paper presents a rigorous evaluation of the *precision* of an automated signal mining pipeline specifically engineered to detect a broad spectrum of operational issues and user experience friction. Our primary objective is to quantify how reliably the signals generated by this pipeline correspond to actual, verifiable problems within the system. We focus on precision because it directly measures the trustworthiness and actionability of the pipeline’s output, which is paramount for its practical deployment in high-stakes operational environments. High precision ensures that engineering or support teams investigating a signal are doing so for a genuine issue, maximizing their efficiency and fostering confidence in the automated system.

To achieve this evaluation, we employed a systematic methodology centered on human expert validation. We meticulously assessed 20 distinct problem-class signals, purposefully sampled from a larger dataset of 121 system rollups, ensuring coverage of various aspects of system operation and user interaction. An expert (Claude) meticulously reviewed each sampled signal, leveraging extensive domain knowledge to determine its validity (i.e., whether it represented a genuine issue), categorize its nature (e.g., software defect, user experience friction), and assign a

severity level. This human-in-the-loop validation process established the ground truth against which the pipeline’s performance was measured. Our analysis not only provides an overall precision score but also delves into class-specific performance, scrutinizing the pipeline’s efficacy in identifying different types of issues and systematically investigating the root causes of false positives.

The insights derived from this evaluation are intended to offer directional guidance into the pipeline’s current efficacy, underscoring its capabilities in specific problem domains, such as defect identification, and outlining clear areas for targeted improvements. Ultimately, this work aims to enhance the pipeline’s accuracy and utility, thereby strengthening its role in maintaining system health and user satisfaction in real-world operational environments.

II. METHODS

A. Overview of Evaluation Methodology

The primary objective of this study was to rigorously evaluate the precision of an automated signal mining pipeline designed to detect operational issues and user experience friction. As highlighted in the introduction, the true utility of such a pipeline hinges on its ability to accurately identify genuine problems without generating excessive false positives. To establish the ground truth necessary for this evaluation, we employed a systematic human-in-the-loop validation approach. A domain expert meticulously reviewed a representative sample of signals generated by the pipeline, classifying each as either a genuine issue or a false positive, categorizing its nature, and assessing its severity. This expert judgment formed the basis for calculating the pipeline’s precision and conducting a detailed analysis of its performance across different issue types.

B. Automated Signal Mining Pipeline

The automated signal mining pipeline under evaluation is engineered to process vast and diverse operational datasets, including system logs, performance metrics, and user interaction data, to automatically identify patterns indicative of underlying operational problems and user experience friction. As discussed, traditional rule-based monitoring often falls short in complex, dynamic systems. This pipeline leverages advanced analytical techniques to generate “problem-class signals,” which are alerts or indicators suggesting a potential issue requiring investigation. While the internal architecture and specific algorithms employed by the pipeline are beyond the scope of this precision evaluation, its functional output—the generation of these actionable signals—was the direct subject of our assessment. Each signal generated by the pipeline represents a hypothesis that a genuine operational issue or UX friction event has occurred.

C. Dataset and Signal Generation

The signals evaluated in this study were derived from a larger dataset encompassing 121 system rollups. A “system rollup” refers to a comprehensive aggregation of operational data pertaining to a specific time period or deployment cycle of the software system, capturing a wide array of events, metrics, and user interactions. From this aggregated data, the automated signal mining pipeline continuously generated problem-class signals throughout the operational period. These signals varied in type, ranging from potential software defects and infrastructure anomalies to indicators of user experience friction, reflecting the pipeline’s broad detection capabilities. The initial pool of potential signals for evaluation was thus generated autonomously by the pipeline from this extensive real-world operational data.

D. Signal Sampling Strategy

From the complete set of problem-class signals identified by the pipeline across the 121 system rollups, a targeted sample of 20 distinct signals was purposefully selected for in-depth human expert validation. This sampling was not random; rather, it was designed to ensure representativeness across different potential issue categories (e.g., defects, UX friction) and to include signals that the pipeline had flagged with varying degrees of confidence, where such confidence scores were available. The goal was to obtain a diverse cross-section of the pipeline’s output to provide a comprehensive and directional evaluation of its efficacy, allowing for class-specific analysis as outlined in the abstract. Each sampled signal was presented for expert review along with all available contextual data that the pipeline would typically provide to an operator.

E. Human Expert Validation

The core of our evaluation methodology involved human expert validation to establish the ground truth for each sampled signal. This process was critical for accurately determining whether a pipeline-generated signal corresponded to a genuine issue or constituted a false positive.

1. Ground Truth Establishment

An expert, referred to as Claude, meticulously reviewed each of the 20 sampled problem-class signals. Claude possessed extensive domain knowledge regarding the system’s architecture, operational characteristics, and user behavior. For each signal, the expert was provided with the signal itself, along with all associated metadata and contextual information that the automated pipeline would typically present to an engineer or operator. This included relevant log snippets, performance charts, user interaction sequences, and any other data points that contributed to the signal’s generation. Based on this comprehensive information and leveraging their deep understanding, Claude independently determined the validity of each signal. A signal was deemed a "True Positive" if it corresponded to a verifiable, genuine operational issue or instance of user experience friction. Conversely, a signal was classified as a "False Positive" if, upon expert review, it was determined not to represent a genuine problem or actionable concern. This expert judgment established the definitive ground truth against which the pipeline’s precision was measured.

2. Issue Categorization and Severity Assessment

In addition to determining validity, the expert also categorized each confirmed genuine issue according to its nature. The primary categories included "software defect" (e.g., bugs, crashes, incorrect functionality) and "user experience friction" (e.g., slow load times, confusing workflows, unexpected behavior impacting usability). This categorization enabled a class-specific analysis of the pipeline’s precision. Furthermore, for all confirmed genuine issues, the expert assigned a severity level. Severity was typically categorized on a scale, such as low, medium, or high, reflecting the impact of the issue on system health, user satisfaction, or business operations. This assessment provided further insight into the types of issues the pipeline successfully identified and their relative importance.

F. Evaluation Metrics

The performance of the automated signal mining pipeline was quantified primarily using the metric of precision, complemented by a detailed analysis of false positives.

1. Precision

Precision was chosen as the primary evaluation metric because it directly measures the trustworthiness and actionability of the pipeline’s output, which is paramount for its practical deployment in high-stakes operational environments. Precision is defined as the proportion of true positive signals among all signals identified by the pipeline, calculated as:

$$\text{Precision} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Positives}}$$

A "True Positive" (TP) was defined as a signal that was correctly identified by the pipeline as an issue and subsequently validated by the human expert as a genuine problem. A "False Positive" (FP) was defined as a signal identified by the pipeline as an issue but determined by the human expert to not be a genuine problem. The overall precision was calculated across all 20 sampled signals. Additionally, class-specific precision was computed for key categories, such as "software defect" and "user experience friction," to understand the pipeline’s performance in identifying different types of issues. This allowed us to discern areas of strength and areas requiring improvement, aligning with our goal of providing directional insights.

2. False Positive Analysis

Beyond the quantitative measure of precision, a qualitative analysis of all identified false positives was conducted. For each false positive signal, the expert documented the probable root cause of the pipeline’s misclassification. This analysis aimed to attribute errors primarily to factors such as insufficient contextual information provided to the pipeline for accurate assessment, or ambiguous signal interpretations where the data could plausibly indicate an issue but was benign in the broader context. This detailed investigation into the nature of false positives is crucial for guiding future enhancements to the signal mining pipeline, helping to refine its algorithms and improve its contextual awareness to reduce alert fatigue and increase operator trust.

III. RESULTS

This section presents the detailed findings from the precision evaluation of the automated signal mining pipeline, based on the human expert validation of 20 sampled problem-class signals. The results quantify the pipeline’s overall accuracy, analyze its performance across different issue categories, and provide insights into the nature of both true and false positive detections.

A. Overall Precision of the Pipeline

The evaluation encompassed 20 problem-class signals sampled from a dataset of 121 system rollups, as described in the Methods section. Out of these 20 sampled signals, the human expert (Claude) confirmed 15 as genuine operational issues or instances of user experience friction, while 5 were identified as false positives. Based on these figures, the automated signal mining pipeline achieved an overall precision of 75.0%. This indicates that three out of every four signals generated by the pipeline correspond to a verifiable, actionable problem within the system. This level of precision is crucial for fostering trust in the automated system and ensuring that engineering or support teams can efficiently investigate alerts without being overwhelmed by noise, aligning with the primary objective outlined in the introduction.

B. Class-Specific Precision Analysis

A more granular analysis of the pipeline’s performance reveals significant differences in precision across distinct problem categories, namely software defects and user experience (UX) friction.

1. Defect Class Precision

For signals classified by the pipeline as potential software defects, the evaluation included 4 samples. All 4 of these signals were confirmed by the expert as genuine software defects. This resulted in a remarkable 100% precision for the defect class. This finding strongly suggests that the pipeline is highly robust and reliable in identifying traditional software bugs, crashes, or incorrect functionalities. The algorithms employed by the pipeline appear to be particularly effective in correlating data patterns indicative of defects with actual system anomalies, making defect signals highly trustworthy for immediate investigation and resolution.

2. User Experience Friction Class Precision

In contrast, the pipeline’s performance in identifying user experience friction signals showed a lower, though still substantial, precision. Out of 16 sampled signals categorized as UX friction, 11 were confirmed as genuine issues by the expert, while 5 were determined to be false positives. This yields a precision of 69% for the UX friction class. The higher rate of false positives in this category suggests that while the pipeline successfully identifies many instances of user experience friction, it also generates a notable amount of noise. This indicates a greater challenge in distinguishing genuine UX problems from benign or ambiguous user behaviors, possibly due to the inherent subjectivity and contextual dependency of user experience issues, as discussed in the introduction.

C. Severity Distribution of Confirmed Issues

Of the 15 confirmed genuine issues, their severity levels were categorized by the human expert. The distribution is as follows:

- **High Severity:** 1 issue (6.7%)
- **Medium Severity:** 12 issues (80.0%)
- **Low Severity:** 2 issues (13.3%)

The vast majority of confirmed issues (80.0%) were classified as medium severity. This suggests that the pipeline is effective in identifying issues that, while not immediately critical, are significant enough to warrant attention and can impact system health or user satisfaction over time. The detection of a high-severity issue underscores the pipeline’s capability to flag critical problems. This distribution indicates that the pipeline successfully identifies a broad spectrum of actionable issues, with a strong emphasis on those that require proactive management.

D. Expert Confidence in Judgments

The human expert (Claude) reported high confidence in their judgments for the majority of the sampled signals. Specifically, confidence levels were distributed as follows: 14 judgments were made with high confidence, and 6 with medium confidence. No judgments were made with low confidence. This high level of expert confidence reinforces the reliability of the ground truth established for this evaluation, thereby strengthening the validity of the precision metrics derived.

E. False Positive Analysis

A detailed qualitative analysis of the 5 false positive signals provided critical insights into the limitations of the current pipeline and areas for improvement. These false positives primarily stemmed from two common characteristics:

1. *Missing Context*

Three of the five false positives were directly attributed to missing session context. In these instances, the contextual information retrieved by the pipeline for evaluation was insufficient for the expert to validate the signal as a genuine issue. The pipeline might have detected an anomaly, but without the broader operational or user interaction context, it was impossible to discern whether the flagged behavior was truly problematic or part of a normal, albeit unusual, system state. This highlights the critical importance of comprehensive context retrieval for accurate signal interpretation.

2. *Ambiguous Behavior*

The remaining two false positives were due to ambiguous behavior. The flagged patterns, while potentially indicative of an issue, could also be interpreted as normal or expected behavior depending on subtle nuances of user intent or system design. An illustrative example of this ambiguity was the signal: “Assistant repeatedly asks for contact channel and tone details.” The expert’s reasoning for judging this as a false positive was that “Asking for contact channel and tone details when crafting outreach is standard practice, not a problem.” This demonstrates the challenge in distinguishing between a system redundantly asking for already provided information (a true problem) versus appropriately prompting for necessary details (normal operation). Such cases underscore the need for the pipeline to develop a more sophisticated understanding of expected user workflows and system interactions.

F. True Positive Examples

To further illustrate the pipeline’s capabilities, two examples of confirmed true positive signals are provided:

- **High Severity Example:** The signal “User repeatedly receives ‘Reminder needs attention’ messages” was confirmed as a REAL ISSUE with HIGH severity. The expert reasoned that “A user repeatedly receiving attention messages without resolution path indicates a systematic failure.” This highlights the pipeline’s ability to identify critical systemic failures that directly impact user workflow and satisfaction.
- **Medium Severity Example:** The signal “System repeatedly reports ‘Cron (error)’ but continues functioning” was confirmed as a REAL ISSUE with MEDIUM severity. The reasoning was that “Error logging during normal operations suggests misclassified successes as failures.” This exemplifies the pipeline’s capacity to detect subtle operational inefficiencies or misconfigurations that, while not causing immediate outages, can lead to alert fatigue, obfuscate genuine problems, or indicate underlying instability.

G. Statistical Considerations

Given the sample size of $n = 20$ signals, the results should be interpreted as directional. A 95% Wilson confidence interval for the overall precision of 75.0% is approximately [53%, 89%]. While the sample size provides a valuable initial assessment, a larger sample would yield a narrower confidence interval and more statistically robust conclusions. Nonetheless, the clear difference in per-class precision (100% for defect detection versus 69% for UX friction) provides strong directional evidence that the pipeline’s defect detection mechanisms are highly reliable, whereas its UX friction detection components require further refinement to reduce false positives.

In summary, the automated signal mining pipeline demonstrates a strong overall precision of 75%, indicating its utility in identifying genuine operational issues. Its performance in detecting software defects is exceptionally high (100% precision), positioning it as a reliable tool for bug identification. However, the pipeline exhibits a higher false positive rate when identifying user experience friction (69% precision), primarily due to limitations in contextual information retrieval and the ambiguity inherent in certain behavioral patterns. The majority of confirmed issues are of medium severity, suggesting the pipeline effectively targets actionable problems. These findings provide a clear roadmap for future enhancements, particularly in refining UX friction detection and improving contextual data integration.

IV. CONCLUSIONS

A. Problem statement

The increasing complexity and scale of modern software systems generate an overwhelming volume of operational data, making manual identification of operational issues and user experience (UX) friction intractable. Automated signal mining pipelines offer a promising solution by leveraging advanced analytics to detect problems proactively. However, the practical utility of such systems is critically dependent on their ability to accurately identify genuine issues without generating an excessive number of false positives, which can lead to alert fatigue and erode user trust. This paper addressed the crucial need to rigorously evaluate the precision of an automated signal mining pipeline, quantifying how reliably its generated signals correspond to actual, verifiable problems within a complex system.

B. Methods and datasets

To evaluate the pipeline’s precision, we employed a systematic human-in-the-loop validation methodology. The automated signal mining pipeline, designed to process diverse operational datasets from 121 system rollups, generated “problem-class signals” indicative of potential issues. From this larger set, a targeted sample of 20 distinct problem-class signals was purposefully selected to ensure representativeness across different issue categories. A domain expert, Claude, meticulously reviewed each sampled signal, along with its associated contextual information, to establish ground truth. This expert validation involved determining the signal’s validity (true positive or false positive), categorizing confirmed issues (e.g., software defect, user experience friction), and assigning a severity level. The primary evaluation metric was precision, calculated as the proportion of true positives among all detected signals. Class-specific precision was also computed, complemented by a detailed qualitative analysis of false positives to identify their root causes.

C. Results obtained

The evaluation revealed that the automated signal mining pipeline achieved an overall precision of 75.0%, indicating that three out of four detected signals corresponded to genuine operational issues or instances of user experience friction. A class-specific analysis highlighted significant differences in performance: the pipeline demonstrated an exceptionally robust 100% precision for the detection of software defects, confirming all 4 sampled defect-class signals as genuine. In contrast, signals related to user experience friction showed a precision of 69%, with 11 out of 16 sampled signals confirmed as genuine issues and 5 identified as false positives. The majority of confirmed issues (80.0%) were classified as medium severity, with 6.7% high severity and 13.3% low severity, suggesting the pipeline effectively targets actionable problems. False positive analysis indicated that these errors primarily stemmed from insufficient contextual information (3 out of 5 false positives) or ambiguous signal interpretations (2 out of 5 false positives), where benign behavior was misclassified. Expert confidence in judgments was predominantly high, reinforcing the reliability of the established ground truth. While the sample size of 20 signals provides directional insights, the statistical confidence interval for overall precision was [53%, 89%].

D. Learnings from the paper

This study offers several key learnings regarding the efficacy and areas for improvement of automated signal mining pipelines. Firstly, the overall precision of 75.0% confirms the pipeline's substantial utility as a tool for automated issue detection, capable of reducing manual effort and enabling proactive system management. Secondly, the outstanding 100% precision in identifying software defects underscores the pipeline's particular strength and reliability in detecting traditional bugs and system malfunctions, making it a highly trustworthy component for defect identification within operational environments. Thirdly, the lower precision for user experience friction signals (69%) highlights a critical area for future enhancement. The analysis of false positives revealed that improving the retrieval and integration of comprehensive contextual information, along with refining algorithms to better differentiate ambiguous but benign user or system behaviors from genuine UX friction, is paramount for increasing accuracy in this category. Finally, the pipeline's effectiveness in identifying a broad spectrum of issues, predominantly of medium severity, confirms its capability to flag actionable problems that contribute to sustained system health and user satisfaction. These findings provide a clear roadmap for targeted improvements, ultimately aiming to enhance the pipeline's overall accuracy, reduce alert fatigue, and solidify its role in maintaining high system resilience and user experience standards.