AUTOGETS: AUTOMATED GENERATION OF TEXT SYNTHETICS FOR IMPROVING TEXT CLASSIFICATION

Anonymous authors

Paper under double-blind review

Abstract

When developing text classification models for real world applications, one major challenge is the difficulty to collect sufficient data for all text classes. In this work, we address this challenge by utilizing large language models (LLMs) to generate synthetic data and using such data to improve the performance of the models without waiting for more real data to be collected and labelled. As an LLM generates different synthetic data in response to different input examples, we formulate an automated workflow, which searches for input examples that lead to more "effective" synthetic data for improving the model concerned. We study three search strategies with an extensive set of experiments, and use experiment results to inform an ensemble algorithm that selects a search strategy according to the characteristics of a class. Our further experiments demonstrate that this ensemble approach is more effective than each individual strategy in our automated workflow for improving classification models using LLMs.

023 024 025

004

010 011

012

013

014

015

016

017

018

019

021

1 INTRODUCTION

026 A critical impediment to developing robust text classification models for real-world applications is 027 the pervasive challenge of class imbalance and data scarcity, particularly for underrepresented text categories. Many industrial applications, such as ticketing systems, require classification models to 029 process large volumes of unstructured text data, such as problem descriptions and user comments, which are often heavily imbalanced in class sizes. In modern industrial environments, ticketing sys-031 tems play a vital role in managing and resolving technical issues, service requests, and operational incidents (Al-Hawari & Barham (2021)). As shown in the workflow (Figure 1), models are initially 033 trained on a set of labeled tickets, but newly introduced or infrequent classes often arise after de-034 ployment, necessitating manual classification and correction. This reliance on manual intervention for new or underrepresented classes creates operational bottlenecks and impairs model adaptation to evolving data distributions. Over time, the model's performance degrades, particularly for small or 036 specialized categories, as obtaining balanced and adequately labeled data across all classes remains 037 challenging. Consequently, models often fail to generalize effectively across diverse data distributions, especially for underrepresented categories Gandla et al. (2024). Traditionally, addressing data scarcity involves collecting additional real-world data, which can be both time-consuming and 040 resource-intensive, especially for rare or newly introduced classes. Furthermore, the manual label-041 ing of such data introduces additional delays and costs Li et al. (2022). In this context, synthetic 042 data generation has emerged as a promising solution to address class imbalance and data scarcity, 043 particularly for underrepresented classes. By augmenting training datasets with synthetic samples, 044 models can achieve improved performance and generalization across different categories.

Synthetic data has gained popularity in recent years as a way to overcome the limitations of real-world data, which can be scarce, sensitive, or expensive to obtain (Patki et al. (2016)). Research in this area consistently highlights the potential of synthetic data to enhance the performance of ML models across diverse fields (Lu et al. (2023)), addressing challenges such as data shortages in computer vision and NLP (Mumuni et al. (2024)), generating diverse datasets in medical imaging (Frid-Adar et al. (2018b)), and providing safe training scenarios for autonomous driving systems (Song et al. (2023)). Its utility extends to financial modeling for algorithm testing under simulated market conditions and cybersecurity for developing threat detection systems (Potluru et al. (2023); Chalé & Bastian (2022)). In the domain of text analysis, synthetic data has been increasingly employed to enhance ML models, particularly in tasks such as text classification, sentiment analysis,

and natural language understanding. Moreover, researchers have shown generating synthetic samples with only targeted data examples could more effectively improve the model (Jin et al. (2024)).
 However, the process of identifying these optimal data examples often requires substantial domain expertise and manual effort, making it time-consuming and less scalable for real-world applications.



Figure 1: The workflow for developing and deploying a classification model in an industrial ticketing system, and the main obstacles impacting on the performance of the model.

To overcome these challenges, this paper introduces *Automated Generation of Text Synthetics* (AutoGeTS), an algorithmic solution that automates the search for optimal data examples based on specific improvement objectives, eliminating the need for human intervention. Through experiments with three search strategies and four objective functions, we identify key patterns between optimal strategy-objective combinations and data characteristics. We propose an ensemble algorithm that effectively improves text classification models across various real-world tasks.

079 080

081

060

061

062

063

064 065

066

067

068

069

071

072 073

074

075

076

077

2 RELATED WORK

082 Synthetic data is increasingly used as a powerful tool for generating realistic datasets to enhance the 083 performance of the task across various domains (Meier et al. (1988); Bersano et al. (1997)). Early 084 synthetic data generation methods include bootstrapping (Efron (1992); Breiman (1996)), which 085 resamples from original data to estimate distributions and reduce variance in predictions, proved effective for a range of predictive algorithms including tree-based models (Sutton (2005)). How-087 ever, bootstrapping couldn't introduce new patterns. The Synthetic Minority Over-sampling Tech-880 nique (SMOTE) (Chawla et al. (2002)) advanced imbalanced dataset handling but risked overfitting. Data augmentation (Jaderberg et al. (2014)) improved model robustness by transforming existing 089 data points, increasing diversity, yet still limited to patterns in the original dataset. The advent of 090 deep learning introduced more sophisticated techniques, notably Generative Adversarial Networks 091 (GANs) by Goodfellow et al. (Goodfellow et al. (2014)), which generate highly realistic synthetic 092 data capturing dataset complexity. Studies have shown models trained on GAN-generated synthetic data often perform comparably to those trained on real data in various predictive tasks (Zhang et al. 094 (2017); Cortés et al. (2020)). Frid-Adar et al. (Frid-Adar et al. (2018a)) enhanced liver lesion diag-095 nosis using GAN-generated images, while Yale et al. (2020) demonstrated comparable performance 096 using GAN-generated synthetic electronic health records for ICU patient predictions.

GANs have been extensively used for synthetic text generation. For instance, Croce et al. (2020) 098 demonstrated their effectiveness in generating realistic text for NLP tasks, while He et al. (2022) explored task-specific text generation. However, GAN-generated data for text classification often 100 lacks semantic coherence and relevance to specific tasks (Torres (2018)). Recent advancements in 101 large language models (LLMs), such as GPT-2 (Croce et al. (2020)), provide new approaches to 102 overcome these limitations. LLMs excel in few-shot and zero-shot learning (Brown (2020); Wang 103 et al. (2021)), adapting to unseen tasks and generating contextually relevant data that improves model 104 robustness. Yoo et al.'s GPT-3Mix (Yoo et al. (2021)) demonstrates LLMs' capability to generate di-105 verse, high-quality synthetic data for text classification through careful prompt engineering. Prompt optimization strategies have shown that carefully crafting input prompts can significantly impact 106 the quality of generated data (Wang et al. (2023)). Automated search techniques for identifying the 107 most effective prompts, such as those used in AutoPrompt (Shin et al. (2020); Xu et al. (2024)),

offer a potential solution for improving synthetic data generation. Beyond prompt engineering, se-lecting appropriate input examples has emerged as a crucial focus. Selecting example data, either with a uniform distribution or human identification through VIS4ML, to form the prompt for LLM to generate synthetics is shown effective (Li et al. (2023); Jin et al. (2024)). Despite these advance-ments, LLM-generated data still struggles to fully capture real-world diversity, especially in highly subjective tasks (Li et al. (2023)). To address this, we propose AutoGeTS, an automated approach that optimizes input example selection for LLM-generated synthetic data. Designed for real-world business requirements, AutoGeTS reduces human intervention while systematically identifying im-pactful examples, enhancing model performance in scenarios of data scarcity and class imbalance.

3 Methods

3.1 AUTOGETS ARCHITECTURE AND WORKFLOW

Figure 2 illustrates the AutoGeTS architecture. After training and evaluating the original model M0, improvement requirements (overall or class-specific) are determined. For a selected class C, visual encoding is applied to the training dataset. The optimal strategy-objective is employed to select example message sets E_s from C, which are then processed through GPT-3.5's API using a zero-shot prompt template, one message per chat (detailed in Appendix B.1). Each example generates multiple synthetic samples through automated parsing and format cleaning of the LLM responses. These samples are appended to the training set for model retraining. The best-performing model in testing, according to the specified goals, is selected for deployment.



Figure 2: The architecture of AutoGeTS and the workflow for training and improving a model.

3.2 OBJECTIVES FOR MODEL OPTIMIZATION

Ticketing systems deployed in specific organizational environments often face different, sometimes conflicting, requirements. Typical business requirements and related performance metrics include:

- R1. The accuracy of every class should be as high as possible and above a certain threshold. One may optimize a model with a performance metric such as class-based *balanced accuracy* or *F1-score* as the objective function, with each threshold value as a constraint.
- R2. The overall classification accuracy of a model should be as high as possible and above a certain threshold because misclassified messages lead to undesirable consequences. One may optimize a model with a global performance metric, such as overall *balanced accuracy* and overall *f1-score*.
- R3. The recall for some specific classes (e.g., important) should be as high as possible and above a certain threshold in order to minimize the delay due to the messages in such a class being sent to other services. Class-based *recall* is the obvious metric for this requirement. Often one may make a balanced judgment by observing Pareto fronts of *recall* in conjunction with another class-based metric (e.g., *balanced accuracy* or *F1-score*).
- These requirements inform the definition of objective functions and constraints for AutoGeTS optimization. However, because the use of LLMs to generate synthetic data to aid ML (i.e., the workflow in Figure 2) is a recent approach, it is necessary to understand how different example selection algorithms for LLMs may impact the optimization.

162 3.3 STRATEGIES FOR EXAMPLE SELECTION

164 Defining the search space for example selection is critical, especially when augmenting datasets with 165 synthetic examples. This space includes all possible subsets of training data $D = x_1, x_2, \ldots, x_n$, 166 each data dot labeled with a specific class. The primary objective is to identify the optimal subset 167 $W^* \subseteq D$ that maximizes performance metrics when used for synthetic data generation via LLM. 168 With $2^n - 1$ possible subsets for *n* examples, exhaustive search becomes intractable, necessitating 169 heuristic strategies.

The general goal can be formulated as an ideal multi-objective optimization problem:

$$W^* = \arg\max_{W \subseteq D} J(W) \tag{1}$$

where J(W) is the objective function measuring the performance of a retrained model M(W) using synthetic data generated from subset W:

$$J(W) = w_1 \cdot \text{Recall}(M(W)) + w_2 \cdot \text{BalancedAccuracy}(M(W)) + w_3 \cdot \text{F1}(M(W))$$
(2)

where w_1, w_2 , and w_3 reflect metrics weights in the overall objective. Given the practical challenges in defining such a compound objective, we employ a simplified, single-metric function:

$$W^* = \arg\max_{W \subset D} J'(W) \tag{3}$$

where J'(W) represents one of the following metrics: Class-based Recall (CR), Class-based Balanced Accuracy (CBA), Overall Balanced Accuracy (OBA), and Overall F1-Score (OF1).

Thus, the policy for selecting optimal subsets involves two core components:

- 1. A strategy for selecting subsets of examples for synthetic data generation and retraining.
- 2. An evaluation metric, J'(W), to be maximized as the objective for the search towards the optimal subset W^* .

The challenge is to efficiently search the space while balancing between computational cost (i.e., cumulative model retraining time) and performance improvement (i.e., maximum gain in J(W)). Initial experiments with random selection yielded limited improvements, particularly for class-based metrics (detailed in Appendix B.1.3). Thus, we explore three primary strategies to optimize the subset selection: brute-force (Sliding Window, SW), gradient-based (Hierarchical Sliding Window, HSW), and evolutionary algorithms (Genetic Algorithm, GA), as illustrated in Figure 3.

200 3.3.1 SLIDING WINDOW (SW)

The Sliding Window (SW) strategy represents a brute-force approach, where the search space is exhaustively segmented into "windows" or subsets. For each window $W_k \subseteq D$, synthetic data is generated, the model is retrained, and the performance is evaluated based on the objective function $J'(W_k)$. The goal is to identify the window W_k^* that yields the maximum improvement:

$$W_k^* = \arg \max_{W_k \subset D} J'(W_k) \tag{4}$$

The brute-force nature of SW ensures that no region of the search space is neglected, but the cost in terms of time and computational resources can become prohibitive.

211 3.3.2 HIERARCHICAL SLIDING WINDOW (HSW)

The Hierarchical Sliding Window (HSW) strategy builds on the principles of hierarchical selection, offering a more computationally efficient approach by incrementally narrowing the search space to promising regions. At each level l, the current search space is partitioned into smaller windows $W_{k,l}$. For each window, synthetic data is generated, the model is retrained, and the performance

205 206 207

170

171 172

173

174

175

181 182

183

187

188 189

190

191

192

216 is evaluated. Only the windows with the highest objective function values are selected for further 217 hierarchical subdivision in the next level: 218

$$W_{k,l+1}^{*} = \arg \max_{W_{k,l+1} \in \text{Subspace}(W_{k,l})} J'(W_{k,l+1})$$
(5)

The process repeats until improvement in J'(W) plateaus or a predefined stopping criterion is met. HSW thus is akin to a targeted optimization approach that progressively homes in on the optimal subset W^* , balancing thorough exploration with reduced computational complexity compared to the brute-force SW method.



Figure 3: The three examples subset selection strategies.

3.3.3 GENETIC ALGORITHM (GA)

The Genetic Algorithm (GA) begins by initializing a population

- 244 245
- 246

248

249

250

240 241

242 243

219 220 221

222

224

225

247

where each candidate solution $S_i \subseteq D$ represents a subset of the training data D, encoded as a priority value-based chromosome. Each above threshold element indicates that the corresponding data example is included in the subset.

 $P = \{S_1, S_2, \dots, S_m\}$

The GA evolves this population over generations, guided by a fitness score F(S) defined as the 251 objective function J'(S), derived from the AutoGeTS process and subsequent performance evalua-252 tion of the retrained model M(S). The algorithm applies three main genetic operators: Selection: 253 At each generation, Lexicase selection Spector (2012) and Clustered Tournament selection Xie & 254 Zhang (2012) are employed to select individuals into the mating pool based on their fitness, where 255 Lexicase selection evaluates the F(S) of input class and the J'(S) of other randomly chosen classes. 256 Crossover: Weight Mapping Crossover Gen et al. (2006) is used to combine two parent solutions 257 S_i and S_j from the mating pool to produce offspring O_k for local exploration. Mutation: Adaptive 258 Polynomia Mutation Si et al. (2011) is applied to offspring O_k to introduce variability for global 259 search. The GA repeats the selection, crossover, and mutation process until reaching a specified number of generations or a convergence criterion. The subset S^* that maximizes the fitness score: 260

261 262

263

264

$$S^* = \arg \max_{S_i \in P} F(S_i) \tag{7}$$

(6)

where $F(S_i) = J'(S_i)$, is finally retrieved. For further details and the step-by-step breakdown of 265 HSW and GA algorithms, refer to Algorithm 1 and 2 in Appendix C.1. 266

Each AutoGeTS run targets a single class C_i , aiming to improve its specific or overall performance 267 through synthetic sample addition. However, class interactions in synthetic data generation have 268 been observed; Jin et al. (2024) found that synthetic data for one class can improve performance for 269 others. Given these interactions and the collective contribution of all classes to overall classification performance, an ensemble algorithm applying AutoGeTS across multiple classes is necessary to optimize both class-specific and overall performance.

- 273
2743.4Ensemble Algorithm
- ²⁷⁵ The ensemble algorithm depends on the specific business requirements, as outlined in Section 3.2:

276 To lift all classes performances above a threshold (R1): Iteratively apply AutoGeTS to each 277 underperforming class (C_{low}) with optimal strategy-objective combination, in the order that most 278 likely improves class performance. In each iteration, append synthetic samples from the optimal 279 retrained model to the training set, and maintain improvements for processed C_{low} above a specified 280 threshold. Terminate when unable to maintain improvements. To improve overall classification 281 accuracy (R2): The same process is applied to each class (C), except the class order that most 282 likely improves overall performance is used, and the algorithm terminates when overall performance 283 plateaus. To improve specific important class's performance (R3): For an important class (IC), identify related classes (RC) that could enhance IC performance with AutoGeTS. Apply AutoGeTS 284 to IC and RC iteratively with optimal strategy-objective combinations, in the order that prioritizes 285 *IC* performance improvement. Terminate when *IC* performance plateaus. 286

Experimental determinations include performance thresholds, *RC* identification, class order, and
optimal strategy-objective combinations, which will be studied in Section 4. The specific details and
the step-by-step breakdown of the algorithms are provided in Algorithm 3, 4, and 5 in Appendix C.2.

290 291

292 293

294

295 296

297 298

299 300

301

302

303

305

306

307

308

310

311

312

4 EXPERIMENTS AND RESULTS

AutoGeTS is evaluated through 1 GPU hour fixed-time experiments to improve M0 in meeting business requirements and to determine optimal strategy-objective policy as outlined in Section 3.2.

4.1 EXPERIMENT SETUP

Class Class Size Balanced Accuracy F1-Score Recall T211350 0.950 0.941 0.921 8529 0.986 0.979 0.977 T1 T3 4719 0.952 0.914 0.922 T5 2755 0.889 0.794 0.794 T7 1963 0.883 0.780 0.766 T6 1888 0.821 0.623 0.665 T10 1699 0.761 0.540 0.554 Т9 0.747 1466 0.861 0.680 T4 1387 0.899 0.801 0.859 T8 0.828 1028 0.665 0.672 T14 0.772 0.548 764 0.607 0.452 T15 543 0.726 0.596 T11 471 0.973 0.947 0.967 0.484 358 0.742 0.608 T12 T13 180 0.333 0.666 0.469 39100 0.923 Overall 0.856 0.856

Table 1: Original CatBoost Model M0 Performance

313 314

We evaluated the AutoGeTS framework using a dataset from an enterprise IT support ticketing system, comprising 39,100 entries labeled into 15 task classes. The dataset is highly imbalanced, with some classes representing less than 1% of the total entries. To mitigate the effect of this imbalance, we split the dataset into 80% for training/validation and 20% for testing, with a further 80-20 split on the training set for validation. The imbalanced nature of the dataset mirrors real-world challenges faced by classification systems in industrial applications.

We used GPT-3.5 (version: 2023-03-15-preview) to generate synthetic text, employing parameters such as temperature = 0.7, max tokens = 550, top p = 0.5, frequency penalty = 0.3, and presence penalty = 0.0. Comparative experiments with the Easy Data Augmentation (EDA) tool (Wei & Zou (2019)), a traditional data augmentation method, demonstrated that while AutoGeTS improved performance with both approaches, LLM-based workflow yielded superior results (detailed in Appendix E). For the baseline classification model, we utilized CatBoost with fixed hyperparameters (300 iterations, learning rate = 0.2, depth = 8, L2 leaf regularization = 1) to ensure consistency across all retrained models. More detailed M0 analysis and prompt are provided in Appendix A and B.1. The effectiveness of AutoGeTS was evaluated using class-based balanced accuracy, recall, and F1-score for local performance, as well as overall balanced accuracy and F1-score for global performance. The performance of the original CatBoost model M0 is shown in Table 1.

Table 2: Performance Comparison with M0, comparing Overall and Class Balanced Accuracy.

Class	Class	M0	Sliding	Window	Hierarch	ical SW	Genetic Algorithm		
Name	Size	Bal Acc	Overall	Class	Overall	Class	Overall	Class	
T2	11350	0.950	▲0.0030	▲0.0050	▲0.0034	▲0.0048	▲0.0009	▼0.0010	
T1	8529	0.986	▲0.0028	▲0.0005	▲0.0029	▲0.0005	▲0.0018	▼0.0018	
T3	4719	0.952	▲0.0030	▲0.0058	▲0.0027	▲0.0062	▲0.0029	▲0.0069	
T5	2755	0.889	▲0.0032	▲0.0189	▲0.0034	▲0.0140	▲0.0010	▲0.0059	
T7	1963	0.883	▲0.0036	▲0.0228	▲0.0034	▲0.0226	▲0.0012	▼0.0026	
T6	1888	0.821	▲0.0035	▲0.0190	▲0.0030	▲0.0196	▲0.0015	▲0.0073	
T10	1699	0.761	▲0.0034	▲0.0281	▲0.0044	▲0.0247	▲0.0027	▲0.0208	
T9	1466	0.861	▲0.0036	▲0.0147	▲0.0027	▲0.0191	▲0.0026	▲0.0077	
T4	1387	0.899	▲0.0029	▲0.0304	▲0.0033	▲0.0369	▲0.0036	▲0.0323	
T8	1028	0.828	▲0.0030	▲0.0321	▲0.0029	▲0.0358	▲0.0020	▲0.0142	
T14	764	0.772	▲0.0023	▲0.0326	▲0.0029	▲0.0395	▲0.0019	▲0.0396	
T15	543	0.726	▲0.0033	▲0.0456	▲0.0034	▲0.0446	▲0.0037	▲0.0533	
T11	471	0.973	▲0.0030	▲0.0054	▲0.0030	▲0.0053	▲0.0039	▲0.0053	
T12	358	0.742	▲0.0037	▲0.0699	▲0.0032	▲0.0772	▲0.0036	▲0.0775	
T13	180	0.666	▲0.0030	▲0.0443	▲0.0037	▲0.0548	▲0.0034	▲0.0548	

4.2 Performance Improvements Overview

The AutoGeTS framework yielded significant improvements in both local and global performance
 metrics, effectively addressing the class imbalance problem evident in Table 2.

351 Smaller, underrepresented classes experienced the largest improvements. For instance, T13's balanced accuracy increased by 5.48 percentage points (pp) from 66.6% in M0, while T12 showed a 352 7.75 pp improvement from 74.2%. In contrast, larger classes like T1 (98.6%) and T2 (95%) saw 353 only marginal gains of around 0.3 pp. This demonstrates AutoGeTS's ability to significantly im-354 prove underrepresented classes without affecting the performance of well-represented ones. This 355 balance is crucial in maintaining overall system performance. The overall balanced accuracy im-356 proved consistently and comparably among classes, with T10 showing the highest overall balanced 357 accuracy improvement of 0.44 pp. This underscores AutoGeTS's synergistic effect, where class-358 based improvements translate to overall performance gains, with minimal trade-offs between local 359 and global performance improvements (see Appendix D.1).

360 361 362

331

4.3 COMPARISON OF EXAMPLE SELECTION STRATEGIES AND OPTIMIZATION OBJECTIVES

Figure 4 compared the performance of three search strategies—SW, HSW, GA—and four objectives—maximizing CR, CBA, OBA, OF1—with respect to both local (class-specific) and global (overall) metrics. Each bar chart is divided into 4 sections for 4 performance metrics, with the four bars each representing the maximum improvement in the objective of maximizing CR, CBA, OBA, or OF1. The choice of strategy-objective played a critical role in the effectiveness of AutoGeTS, with each demonstrating distinct advantages depending on the size of the target class, and their performance trajectories over retraining time (see Appendix D.2).

For larger classes, HSW consistently yielded the best results, such as the highest class T1 balanced
 accuracy improvement of 0.5%. HSW's progressive narrowing of the search space proves effective
 for larger data sets where an exhaustive search is computationally prohibitive. Objective-wise, max imizing CR or CBA each best improved its respective metric, while maximizing OBA or OF1 both
 led to the best improvements in global metrics.

GA strategy proved superior for smaller and mid-sized classes, as shown by T13 and T12's highest
 balanced accuracy gains. GA's evolutionary nature generates diverse synthetic samples, crucial for
 small data sets. For these classes, maximizing CBA outperformed other objectives in local metrics,
 while OBA or OF1 maximization equally improved global metrics, except for T11, T12, and T15.

SW showed moderate performance across mid-sized classes, such as improving T5 and T6 balanced accuracy by up to 1.89% and 1.90%, respectively. SW offers a balanced trade-off between computational cost and performance improvement for mid-sized data sets. For these classes, maximizing CR or CBA equally improved both local metrics, and the same applies to maximizing OBA or OF1.



Figure 4: Comparison of improvements across 4 metrics for all classes, showing best-performing strategies (SW, HSW, GA) and highest improvement values for each objective

4.4 PARETO ANALYSIS FOR REPRESENTATIVE CLASSES



Figure 5: Improvements and Pareto Fronts in Class Recall vs Class Balanced Accuracy for topics T1, T6, and T13, showing models maximizing TR and TBA (note different axis ranges).

Figure 5 shows trade-offs between Class Recall (CR) and Class Balanced Accuracy (CBA) improvements for classes T1 (large), T6 (mid-sized), and T13 (small) when maximizing either CR or CBA.

CR improvements generally correlate with CBA gains, varying by class size. Large class T1 shows
 the largest divergence (CR +3.5%, CBA +1.2%). The more synthetic samples needed to impact
 a large class directly affects recall but indirectly specificity, leading to larger CR and CBA divergence. Mid-sized class T6 demonstrates aligned improvements, with targeting CR increasing CBA

by 2.9%, indicating a less critical objective choice between maximizing CR or CBA. Small class T13 exhibits substantial improvements in both CR and CBA regardless of which metric is maximized, as maximizing CR improved recall by 33.3% and CBA by 8.2%, reflecting the effectiveness of diverse synthetic samples and the GA strategy for underrepresented classes and small, imbalanced datasets. Practically, maximizing CR is preferable for large important classes to minimize misclassification delay, while either CR or CBA optimization can be effective for middle and small classes.

ENSEMBLE ALGORITHM AND FURTHER EXPERIMENTATION

5.1 SUMMARY OF STRATEGY-OBJECTIVE COMBINATIONS

Building on our analysis in Section 4.3, we synthesize the effectiveness of different strategyobjective combinations across varying class sizes and performance metrics.

Table 3: Optimal Strategy-Objective Combinations across Classes

	Cla	ass Performance	Overall Perform	nance			
	Topic Recall	Topic Balanced Accuracy	Overall Balanced Accuracy Overall F1-Score				
T2	USW TD	HSW-TBA/OBA	HSW-TBA/OBA				
T1	115 w-1K	HSW-TBA	HSW-OBA/C	DF1			
T3		GA-TR/TBA	SW-OBA/O	F1			
T5		SW-TR/TBA	HSW-OBA/C)F1			
T7		SW-TBA					
T6	SW-TR	HSW-OBA/OF1	5W-0DA/011				
T10		SW-TR/TBA	HSW-OBA/OF1				
T9			SW-OBA/OF1				
T4	H	HSW-TR/TBA	GA-OBA/OF1				
T8			SW-OBA/OF1				
T14			HSW-OBA/C	DF1			
T15			GA-OF1				
T11		GA-TBA	GA-OBA				
T12			GA-OF1				
T13			HSW-OBA/OF1				

Table 3 summarizes optimal Strategy-Objective combinations for improving 4 metrics across all classes, serving as a look-up table for the ensemble algorithm in Section 5.

5.2 LOCAL AND GLOBAL METRICS IMPROVEMENT ACROSS CLASSES

We now examine the broader impacts of applying AutoGeTS to individual classes, considering interclass effects and overall system performance. Figure 6 illustrates inter-class interaction and overall effects, guiding class order determination for three requirements: R1 (improve each class) orders classes by descending diagonal elements of class balanced accuracy improvements; R2 (improve overall performance) sorts classes by their overall balanced accuracy improvement (Figure 6b); and R3 (improve an important class) orders related classes RC by their improvement on the class bal-anced accuracy of important class IC, according to the IC's column of Figure 6a. The determined class orders for the three requirements are detailed in Appendix C.2.

ENSEMBLE CASE STUDY RESULTS FOR R3 5.3

We applied ensemble AutoGeTS to improve important class T13 following Algorithm 5 in Ap-pendix C.2, with related classes T12, T10, T11, and T5 identified and ordered based on T13's column of Figure 6a. We selected optimal strategy-objective combinations for each class from Ta-ble 3. Benchmarks included iterating single strategies and random combinations. We conducted three runs with different train-validation splits and random seeds. T13 balanced accuracy was used as the performance metric instead of T13 recall, as they align well for T13 while reflecting changes in other classes (negative instances for T13) not captured by the recall.

The Ensemble Algorithm in Figure 7 achieved the highest T13 balanced accuracy and second-highest global improvements, with faster and consistently higher T13 performance gains. While





Figure 7: Average maximum cumulative T13 and overall improvements for 5 ensemble sequences.

CONCLUSIONS

513 514

515 516

517

518 519

520

523

524 This work introduces AutoGeTS, an automated framework optimizing example data selection for 525 synthetic text generation using Large Language Models (LLMs). AutoGeTS significantly enhances 526 the level of automation, reducing human efforts in selecting effective examples. This approach addresses class imbalance and data scarcity challenges in real-world text classification tasks. Exper-527 iments demonstrate AutoGeTS's effectiveness in improving both local and global performance met-528 rics. Using Sliding Window, Hierarchical Sliding Window, and Genetic Algorithm, significant im-529 provements in inadequately-sampled classes are observed without compromising well-represented 530 ones. The ensemble algorithm for selecting the most suitable strategies according to the results of 531 earlier iterations facilitates more efficient optimization processes in later iterations. This approach 532 that treats earlier testing results as useful knowledge in optimization will likely have a wider appli-533 cation in ML model deployment. These findings establish AutoGeTS as an effective solution for 534 enriching training data with synthetic samples to meet real-world requirements for the performance 535 of ML models with limitations in data collection. This work confirms that the automated approach 536 can perform better than human-centric processes in terms of both effectiveness and efficiency. Mean-537 while, there is a need to conduct more large-scale experiments and analyze the experiment results in order to understand and explain how different synthetic data sets improve or undermine an ML 538 model. Future work may also explore multi-objective optimization strategies, more advanced ensemble algorithms, and applications of this approach in other ML domains.

540 REFERENCES 541

547

550

556

558

559

563

564

565

569

571

576

577

- Feras Al-Hawari and Hala Barham. A machine learning based help desk system for it service 542 management. Journal of King Saud University-Computer and Information Sciences, 33(6):702-543 718, 2021. 544
- Tom Bersano, Brad Clement, and Leonid Shilkrot. Synthetic data for testing in databases. University 546 of Michigan, 1997.
- Leo Breiman. Bagging predictors. *Machine learning*, 24:123–140, 1996. 548
- 549 Tom B Brown. Language models are few-shot learners. arXiv preprint arXiv:2005.14165, 2020.
- Marc Chalé and Nathaniel D Bastian. Generating realistic cyber data for training and evaluating 551 machine learning classifiers for network intrusion detection systems. Expert Systems with Appli-552 cations, 207:117936, 2022. 553
- 554 Nitesh V Chawla, Kevin W Bowyer, Lawrence O Hall, and W Philip Kegelmeyer. Smote: synthetic 555 minority over-sampling technique. Journal of artificial intelligence research, 16:321–357, 2002.
 - Andoni Cortés, Clemente Rodríguez, Gorka Vélez, Javier Barandiarán, and Marcos Nieto. Analysis of classifier training on synthetic data for cross-domain datasets. IEEE Transactions on Intelligent Transportation Systems, 23(1):190–199, 2020.
- Danilo Croce, Giuseppe Castellucci, and Roberto Basili. Gan-bert: Generative adversarial learning for robust text classification with a bunch of labeled examples. In Proceedings of the 58th annual 561 meeting of the association for computational linguistics, pp. 2114–2119, 2020. 562
 - Bradley Efron. Bootstrap methods: another look at the jackknife. In Breakthroughs in statistics: Methodology and distribution, pp. 569–593. Springer, 1992.
- Maayan Frid-Adar, Idit Diamant, Eyal Klang, Michal Amitai, Jacob Goldberger, and Havit 566 Greenspan. Gan-based synthetic medical image augmentation for increased cnn performance 567 in liver lesion classification. Neurocomputing, 321:321-331, 2018a. 568
- Maayan Frid-Adar, Eyal Klang, Michal Amitai, Jacob Goldberger, and Hayit Greenspan. Synthetic 570 data augmentation using gan for improved liver lesion classification. In 2018 IEEE 15th international symposium on biomedical imaging (ISBI 2018), pp. 289–293. IEEE, 2018b.
- 572 Phani Krishna Kollapur Gandla, Rajesh Kumar Verma, Chhabi Rani Panigrahi, and Bibudhendu 573 Pati. Ticket Classification Using Machine Learning, pp. 487–501. Springer Nature Singapore, 574 2024. ISBN 9789819950157. 575
 - Mitsuo Gen, Fulya Altiparmak, and Lin Lin. A genetic algorithm for two-stage transportation problem using priority-based encoding. OR spectrum, 28:337–354, 2006.
- 578 Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, 579 Aaron Courville, and Yoshua Bengio. Generative adversarial nets. Advances in neural information 580 processing systems, 27, 2014.
- Xuanli He, Islam Nassar, Jamie Kiros, Gholamreza Haffari, and Mohammad Norouzi. Generate, 582 annotate, and learn: Nlp with synthetic text. Transactions of the Association for Computational 583 Linguistics, 10:826-842, 2022. 584
- 585 M Jaderberg, K Simonyan, A Vedaldi, and A Zisserman. Synthetic data and artificial neural net-586 works for natural scene text recognition. In NIPS Deep Learning Workshop. Neural Information Processing Systems, 2014.
- 588 Siwei Jiang, Yew-Soon Ong, Jie Zhang, and Liang Feng. Consistencies and contradictions of performance metrics in multiobjective optimization. IEEE transactions on cybernetics, 44(12):2391-590 2404, 2014. 591
- Yuanzhe Jin, Adrian Carrasco-Revilla, and Min Chen. igaiva: Integrated generative ai and visual 592 analytics in a machine learning workflow for text classification. arXiv preprint arXiv:2409.15848, 2024.

594 Qian Li, Hao Peng, Jianxin Li, Congying Xia, Renyu Yang, Lichao Sun, Philip S Yu, and Lifang 595 He. A survey on text classification: From traditional to deep learning. ACM Transactions on 596 Intelligent Systems and Technology (TIST), 13(2):1–41, 2022. 597 Zhuoyan Li, Hangxiao Zhu, Zhuoran Lu, and Ming Yin. Synthetic data generation with large lan-598 guage models for text classification: Potential and limitations. In The 2023 Conference on Empirical Methods in Natural Language Processing, 2023. 600 601 Yingzhou Lu, Minjie Shen, Huazheng Wang, Xiao Wang, Capucine van Rechem, and Wenqi Wei. 602 Machine learning for synthetic data generation: a review. arXiv preprint arXiv:2302.04062, 2023. 603 Alan K Meier, John Busch, and Craig C Conner. Testing the accuracy of a measurement-based 604 building energy model with synthetic data. Energy and buildings, 12(1):77-82, 1988. 605 Alhassan Mumuni, Fuseini Mumuni, and Nana Kobina Gerrar. A survey of synthetic data augmen-607 tation methods in computer vision. arXiv preprint arXiv:2403.10075, 2024. 608 Neha Patki, Roy Wedge, and Kalyan Veeramachaneni. The synthetic data vault. In 2016 IEEE 609 international conference on data science and advanced analytics (DSAA), pp. 399–410. IEEE, 610 2016. 611 612 Vamsi K Potluru, Daniel Borrajo, Andrea Coletta, Niccolò Dalmasso, Yousef El-Laham, Elizabeth 613 Fons, Mohsen Ghassemi, Sriram Gopalakrishnan, Vikesh Gosai, Eleonora Kreačić, et al. Syn-614 thetic data applications in finance. arXiv preprint arXiv:2401.00081, 2023. 615 Taylor Shin, Yasaman Razeghi, Robert L Logan IV, Eric Wallace, and Sameer Singh. Autoprompt: 616 Eliciting knowledge from language models with automatically generated prompts. In Proceed-617 ings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP). 618 Association for Computational Linguistics, 2020. 619 620 Tapas Si, ND Jana, and Jaya Sil. Particle swarm optimization with adaptive polynomial mutation. 621 In 2011 World Congress on Information and Communication Technologies, pp. 143–147. IEEE, 622 2011. 623 Zhihang Song, Zimin He, Xingyu Li, Qiming Ma, Ruibo Ming, Zhiqi Mao, Huaxin Pei, Lihui Peng, 624 Jianming Hu, Danya Yao, et al. Synthetic datasets for autonomous driving: A survey. IEEE 625 Transactions on Intelligent Vehicles, 2023. 626 Lee Spector. Assessment of problem modality by differential performance of lexicase selection 627 in genetic programming: a preliminary report. In Proceedings of the 14th annual conference 628 companion on Genetic and evolutionary computation, pp. 401–408, 2012. 629 630 Clifton D Sutton. Classification and regression trees, bagging, and boosting. Handbook of statistics, 631 24:303-329, 2005. 632 D Garcia Torres. Generation of synthetic data with generative adversarial networks. Unpublished 633 doctoral dissertation). Ph. D. Thesis, Royal Institute of Technology, Stockholm, Sweden, 26, 2018. 634 635 Xinyuan Wang, Chenxi Li, Zhen Wang, Fan Bai, Haotian Luo, Jiayou Zhang, Nebojsa Jojic, Eric 636 Xing, and Zhiting Hu. Promptagent: Strategic planning with language models enables expert-637 level prompt optimization. In The Twelfth International Conference on Learning Representations, 638 2023. 639 Zirui Wang, Adams Wei Yu, Orhan Firat, and Yuan Cao. Towards zero-label language learning. 640 arXiv preprint arXiv:2109.09193, 2021. 641 642 Jason Wei and Kai Zou. Eda: Easy data augmentation techniques for boosting performance on text 643 classification tasks. In Proceedings of the 2019 Conference on Empirical Methods in Natural 644 Language Processing and the 9th International Joint Conference on Natural Language Process-645 ing (EMNLP-IJCNLP), pp. 6382–6388, 2019. 646 Huayang Xie and Mengjie Zhang. Parent selection pressure auto-tuning for tournament selection in 647

genetic programming. IEEE Transactions on Evolutionary Computation, 17(1):1–19, 2012.

648 649 650 651	Ran Xu, Hejie Cui, Yue Yu, Xuan Kan, Wenqi Shi, Yuchen Zhuang, May Dongmei Wang, Wei Jin, Joyce Ho, and Carl Yang. Knowledge-infused prompting: Assessing and advancing clinical text data generation with large language models. In <i>Findings of the Association for Computational Linguistics ACL 2024</i> , pp. 15496–15523, 2024.
652 653 654 655	Andrew Yale, Saloni Dash, Ritik Dutta, Isabelle Guyon, Adrien Pavao, and Kristin P Bennett. Generation and evaluation of privacy preserving synthetic health data. <i>Neurocomputing</i> , 416:244–255, 2020.
656 657 658	Kang Min Yoo, Dongju Park, Jaewook Kang, Sang-Woo Lee, and Woomyoung Park. Gpt3mix: Leveraging large-scale language models for text augmentation. In <i>Findings of the Association for Computational Linguistics: EMNLP 2021</i> , pp. 2225–2239, 2021.
659 660 661 662	Yizhe Zhang, Zhe Gan, Kai Fan, Zhi Chen, Ricardo Henao, Dinghan Shen, and Lawrence Carin. Adversarial feature matching for text generation. In <i>International conference on machine learning</i> , pp. 4006–4015. PMLR, 2017.
663 664 665	Eckart Zitzler and Lothar Thiele. Multiobjective evolutionary algorithms: a comparative case study and the strength pareto approach. <i>IEEE transactions on Evolutionary Computation</i> , 3(4):257–271, 1999.
666 667 668	
669 670 671	
672 673 674	
675 676 677	
678 679 680	
681 682 683	
684 685 686	
687 688 689	
690 691 692	
693 694 695	
696 697 698	
699 700 701	

702 APPENDICES

In the following appendices, we provide further experiment results through visualization plots. The experimental data will be made available on GitHub after the double-blind review process. These appendices include:

- A. **Parameters for ML Training** We report the pilot experiments for selecting parameters that were used in the training of the benchmark model M0 (trained without synthetic data). The relatively optimal parameter set was chosen and used for retraining all other models (trained with both collected data and synthetic data).
- B. **Parameters for Example Search** We report the pilot experiments selecting parameters to be used by different algorithms that search message examples to be used as the inputs to LLMs in order to generate synthetic data.
- C. Algorithms We report the detailed algorithm flowcharts for example data subset selection, determined after pilot experiments, and multi-class ensemble algorithms, determined through experiments in Section 4.
- D. Fixed-Time Experiments We report a set of experiments, where three workflows were allowed to use exactly one hour of GPU time for searching message examples, generating synthetic data, training and testing a model in multiple iterations in order to develop a model to improve the benchmark model M0.
- E. Comparison with Traditional Data Augmentation We provide comparative analysis between EDA-based and LLM-based AutoGeTS workflow, examining both the best improvements (overall and class-specific) and the temporal progression of overall balanced accuracy improvement across all 15 classes. The Easy Data Augmentation (EDA) tool (Wei & Zou (2019)) is a traditional data augmentation method.
 - A PILOT EXPERIMENTS: PARAMETER FOR ML TRAINING

Before developing the AutoGeTS framework, we aimed to improve the original CatBoost model, M0, through parameter tuning. A grid search was conducted with the following parameter ranges:

- learning_rate: [0.01, 0.05, 0.1, 0.2, 0.5]
- depth: [4, 6, 8, 10]
- 12_leaf_reg: [1, 3, 5, 10]
- ⁷⁵⁴Five-fold cross-validation was employed, with overall classification accuracy as the primary criterion for evaluating the model performance.

771

772 773

791 792

793 794

796

797

798

799

A.1 **BENCHMARK MO PARAMETERS EXPERIMENTS** 757 758 Parallel Coordinates Plot for Benchmark M0 Catboost Model T3 T5 T7 T6 T10 T9 T4 T8 T14 T15 T11 T12 T13 ovr 0.916 0.785 0.824 0.766 0.556 0.737 0.862 0.713 0.496 0.512 0.954 0.575 0.250 0.854 759 0.780 0.820 0.750 0.550 860 0.700 0.916 924 0.977 0.240 0.850 . 500 0.920 0.976 0.914 4 0.770 815 000 0 84 500 0.220 0.840 761 ,400 0.400 0.600 450 0.200 0 0.920 0.918-0.975 <u>0.912</u> 0.760 0.400 700 3,000 0.550 0.180 0.820 762 0. 350 0.350 0.916 0.974 .500 0.160 0.810 763 0.450 0.300 0.300 0.914 0 800 0.973 0.910 0.300 140 0.80 0.750 0 912 0 910 0.973 0 795 764 50 0.250 908 0.740 0.790 765 0.300 908 0.9 0.200 0.200 0.785 y/290 0 0.150 0.150 0.150 0.150 0.200 60 0.904 0.150 0.100 0.100 07 0.902 .200 0.100 250 50 0 100 0 100 0.900 (100 68 0.100 768 2.000 740 0.050 0.050 0.100 740 500 0.898 0.050 0.050 0 0.050 0.050 .967 0.902 0.710 0.765 4.000⁴1.000 0.730 0.898 0.010 4.000 1.000 0.727 0.8 50 0.000 0.000 0.770 0.0000 0.000 0.000 0.000 0.000 0.000 0.000 0.000 0.000 0.000 0.000 0. 769 0 450 0.966 0.901 0.706 0.763 770

Figure 8: M0 Parameter Experiments Overview.

Figure 8 presents an overview of all experimented M0 CatBoost model parameters 774 775 in a parallel coordinate plot. Each line represents a unique parameter set and its 776 corresponding classification performance. The first three coordinates depict the 777 experimented parameters: learning rate, tree depth, and L2 leaf regularization. The 778 subsequent 15 coordinates (T1 to T15) represent the recall for each of the 15 classes, 779 while the final coordinate shows the overall classification accuracy, which is the 780 primary performance metric. 781



(a) M0 Top Performing Parameters.

Figure 9: M0 Top Performances and Learning Rate = 0.2.

3 5

0.5

Figure 9a highlights the top-performing parameter sets. These sets consistently use a learning rate of 0.2, tree depths of 6 or 8, and L2 leaf regularization values of 1 or 3. Based on these observations, we further examine the performance of learning_rate = 0.2 and determine the optimal values for depth and L2 leaf regularization.

Figure 9b highlights parameter sets with learning_rate = 0.2. All highlighted sets 801 demonstrate good accuracy, including the three best accuracy scores, confirming 802 0.2 as the optimal learning rate among the experimented values. 803

804 Figures 10a and 10b compare model performances between depth = 6 and 805 depth = 8 with learning rate set to 0.2. While both depth values yield good ac-806 curacy, depth = 8 shows slightly superior results overall. 807

808 Figures 11a and 11b compare model performances between 809 $L2_leaf_regularization = 1$ and $L2_leaf_regularization = 3$ with learning rate





Figure 12: M0 PCA Plots for Small Classes T12, T13, T14, T15.

В **PILOT EXPERIMENTS: PARAMETERS FOR EXAMPLE SEARCH**

880 In developing the AutoGeTS framework described in Section 3.1, we found that parameters for both the LLM's synthetic sample generation and the three example 882 selection strategies significantly influenced AutoGeTS performance. To determine optimal parameter sets and understand their impact, we conducted extensive exper-884 iments on each component. 885

886 Given that our objectives for each retrained model were to maximize both overall 887 accuracy and class-specific recall for the chosen class, we employed the Hypervol-888 ume (HV) indicator (Zitzler & Thiele (1999); Jiang et al. (2014)) to evaluate per-889 890 formance. This indicator allows us to compare results across different parameter 891 configurations by considering both class-based recall and overall accuracy simulta-892 neously. 893

894 We implemented 5-fold cross-validation throughout our experiments. In addition to 895 the HV indicator, we tracked the best accuracy and best class recall across all five 896 folds as supplementary performance metrics. 897

SYNTHETIC DATA GENERATION PARAMETER EXPERIMENTS **B** 1

The synthetic data generation process utilizes GPT-3.5 through its API interface. In preliminary experiments, we also evaluated Llama 3 as an alternative language model for synthetic sample generation, which yielded comparable results.

905 For each original text sample, we invoke a new chat session and employ a zero-906 shot approach without providing additional context. The input prompt template for 907 generating synthetic samples follows this format: 908

909 910

898 899

900 901

902

903

904

875 876 877

878 879

881

883

```
'Generate ' + str(num) + ' lines of the data similar to this format data:
        meta_data['text'].values[i] ' + 'put & at the end of each line'
```

911

912 where 'num' is the number of generated samples, meta_data is the input data, and 913 'i' is the data index number. 914

915 Upon receiving the LLM's response, we implement an automated pipeline for 916 cleaning, parsing, and separating the output into 'num' synthetic samples based 917 on the formatting parameters specified in the prompt template. A verification function inspects each synthetic sample for format quality assurance, checking includes: • Empty samples or null responses • Extraneous empty lines or spaces • Correct placement of separation symbols ('&') If 'm' samples fail these quality checks, the generation process is automatically repeated for the same input text, adjusting the prompt to request only the remaining 'm' samples (i.e., 'Generate ' + str(m) + ' lines...'). We investigated the impact of varying 'num' on AutoGeTS performance through a series of experiments. Parallel Coordinates Plot for Number of Synthetic Samples Parameter Tests (3 Topics Overlay)
 Aug. HV F0 HV F1 HV F2 HV F3 HV F4 HV 0 MaxO1 MaxO2 MaxO3 MaxO4 MaxO 0 MaxT1 MaxT2 MaxT3 MaxT4 MaxT4 MaxT4 Vg HV

 0000 0.0019 0.0020 0.0011 0.0028 0.0024
 0.8614
 0.8566 0.8566 0.8564 0.6545 0.9466 0.9399 0.9426 0.9387 0.9413 0.0019

 0.0018 0.0012 0.0026 0.0024
 0.8514
 0.8564 0.8566 0.8566 0.8566 0.8624 0.8540 0.9000 0
 0.0016
 0.0011
 0.0022
 0.0014
 0.0012
 0.0852
 0.8564
 0.8564
 0.8524
 0.8500
 0.8500
 0.8500
 0.8500
 0.8500
 0.8500
 0.8500
 0.8500
 0.8500
 0.8500
 0.8500
 0.8500
 0.8500
 0.8500
 0.8500
 0.8500
 0.8500
 0.8500
 0.8500
 0.8500
 0.8500
 0.8500
 0.8500
 0.8500
 0.8500
 0.8500
 0.8500
 0.8500
 0.8500
 0.8500
 0.8500
 0.8500
 0.8500
 0.8500
 0.8500
 0.8500
 0.8500
 0.8500
 0.8500
 0.8500
 0.8500
 0.8500
 0.8500
 0.8500
 0.8500
 0.8500
 0.8500
 0.8500
 0.8500
 0.8500
 0.8500
 0.8500
 0.8500
 0.8500
 0.8500
 0.8500
 0.8500
 0.8500
 0.8500
 0.8500
 0.8500
 0.8500
 0.8500
 0.8500
 0.8500
 0.8530
 0.8500
 0.8530
 0.8500
 0.8530
 0.8500
 0.8530
 0.8500
 0.8530
 0.8530
 0.8500< 4.5000

Figure 13 presents a parallel coordinate plot of synthetic text generation parameter experiments. The study utilized the Hierarchical Sliding Window approach, focus-ing on classes T2 (orange, largest class), T9 (blue, median size class), and T13 (green, smallest class). The primary parameter under investigation, "Syn Number," represents the number of synthetic text samples generated for each selected origi-nal text data point. The primary criterion, HV, appears as both the second-left and rightmost coordinates in the plot.

2.5000 0.0008 0.0000 0.0012 0.6000 0.0010 0.0010 0.0000 0.60000 0.6000 0.6000 0.6000 0.6000 0.6000 0.6000 0.6000 0.6000 0.6000 0

2.0000 0.0006 0.0004 0.0008 0.0008 0.0008 0.6008 0.6008 0.6550 0.5500 0.

 1.5000
 0.0004
 0.0003
 0.0004
 0.0003
 0.0004
 0.0002
 0.8546
 0.8546
 0.8526
 0.8526
 0.8526
 0.8526
 0.8526
 0.8526
 0.8526
 0.8526
 0.8526
 0.8526
 0.8526
 0.8526
 0.8526
 0.8526
 0.8526
 0.8526
 0.8526
 0.8526
 0.8526
 0.8526
 0.8526
 0.8526
 0.8526
 0.8526
 0.8526
 0.8526
 0.8526
 0.8526
 0.8526
 0.8526
 0.8526
 0.8526
 0.8526
 0.8526
 0.8526
 0.8526
 0.8526
 0.8526
 0.8526
 0.8526
 0.8526
 0.8526
 0.8526
 0.8526
 0.8526
 0.8526
 0.8526
 0.8526
 0.8526
 0.8526
 0.8526
 0.8526
 0.8526
 0.8526
 0.8526
 0.8526
 0.8526
 0.8526
 0.8526
 0.8526
 0.8526
 0.8526
 0.8526
 0.8526
 0.8526
 0.8526
 0.8526
 0.8526
 0.8526
 0.8526
 0.8526
 0.8526
 0.8526
 0.8526
 0.8526
 0.8526<

ume, MaxO: Max Overall Accuracy Improvement, MaxT: Max Topic Recall Improvement).

Figure 13: Synthetic Data Generation Parameter Experiments Overview (F0_HV: Fold 0 Hypervol-

The results indicate that Syn Number = 5 consistently yielded the best Hypervol-ume for all three classes among the tested values. Consequently, we adopted the generation of five synthetic samples per selected original data point for all subse-quent experiments.

B.1.1 PCA PROJECTION OF SYNTHETIC SAMPLES

To verify the effectiveness of generated synthetic samples in addressing class im-balance and data scarcity, we projected these data using the same fitted vectorizer and PCA model used for the original data.

Figure 14 presents updated PCA plots for small classes T12, T13, T14, and T15. The inclusion of synthetic samples significantly increased the number of colored



Figure 14: M0 PCA Plots with Synthetic Samples for Small Classes T12, T13, T14, T15.

data points. Moreover, these additional points appear more densely clustered, suggesting improvement on the previous data scarcity issues.

To illustrate the nature of the synthetic data generated, we present an example of an original data example and its corresponding LLM-generated synthetic samples. The following listing shows the original Spanish text followed by five synthetic samples, demonstrating how the LLM maintains the context and structure while introducing variations in content:

Listing 1: Original data example and LLM-generated synthetic samples

996 Original Data Example: 997 phone_nmb En el ticket REQ0026231 se le solicita acceso a unidades de red pero 998 sique sin poder acceder y figura como resuelto, revisar por favor, la 999 usuaria lleva 2 meses con este problema 1000 Synthetic Samples: 1002 1. El cliente reporta que su cuenta de correo electronico ha sido bloqueada, por favor revisar el caso REQ0027456. 1003 1004 2. Se solicita la instalacion de un software especifico en el equipo del 1005 usuario, el ticket es REQ0028745. 1006 1007 3. El usuario indica que no puede imprimir desde su equipo, se necesita revision del caso REQ0029367. 1008 1009 4. La usuaria reporta problemas con su conexion a internet, el ticket es 1010 REQ0030172. 1011 1012 5. Se requiere la asignacion de permisos adicionales en el sistema para el usuario, el caso es REQ0031298.

1013 1014 1015

983 984 985

986

987 988

989

990

991

992

993 994

995

As evident from Listing 1, the synthetic samples maintain the overall structure of a ticketing system entry while diversifying the reported issues, demonstrating the LLM's ability to generate contextually relevant and varied data.

1018 1019

1016

1017

1020 B.1.2 SYNTHETIC SAMPLES PERFORMANCES WITHOUT EXAMPLES SELECTION 1021

To evaluate the potential of our generated synthetic samples in improving the text classification model M0, we appended all generated synthetic samples to the training set and retrained the CatBoost model M0 for small classes T12, T13, T14, and T15. Tables 4 and 5 present the results of this experiment. We observed that class-based performance often improved when synthetic samples for that class were appended, demonstrating the potential of synthetic data. However, only T12 showed improvement in overall performance. Notably, T13 failed to improve even its class-specific performance.

Table 4: Performance of Retrained Models with T12 and T13 Synthetic Samples

Class	Δ Balance	d Accuracy	ΔR	ecall	Δ F1-Score		
Class	T12	T13	T12	T13	T12	T13	
T1	▼0.0014	▼0.0011	▼0.0034	▼0.0023	▼0.0006	▼0.0009	
T2	▲0.0036	▼0.0033	▲0.0044	▼0.0058	▲0.0053	▼0.0040	
T3	▲0.0015	▼0.0003	▲0.0032	▲0.0011	▲0.0013	▼0.0054	
T4	▲0.0104	▼0.0001	▲0.0214	0.0000	▲0.0064	▼0.0016	
T5	▼0.0015	▼0.0018	▼0.0038	▼0.0038	▲0.0022	▼0.0015	
T6	▼0.0021	▲0.0036	▼0.0027	▲0.0080	▼0.0102	▼0.0003	
T7	▲0.0023	▼0.0039	▲0.0025	▼0.0076	▲0.0161	▼0.0065	
T8	▲0.0046	▼0.0104	▲0.0085	▼0.0212	▲0.0144	▼0.0101	
T9	▼0.0011	▼0.0066	▼0.0036	▼0.0144	▲0.0093	▲0.0014	
T10	▲0.0094	▲0.0043	▲0.0179	▲0.0090	▲0.0192	▲0.0041	
T11	▲0.0001	0.0000	0.0000	0.0000	▲0.0053	0.0000	
T12	▲0.0376	▼0.0234	▲0.0781	▼0.0469	▼0.0497	▼0.0364	
T13	▼0.0333	▼0.0127	▼0.0667	▼0.0222	▼0.0753	▼0.1506	
T14	▲0.0110	▲0.0040	▲0.0222	▲0.0074	▲0.0144	▲0.0184	
T15	▼0.0045	▼0.0298	▼0.0085	▼0.0598	▼0.0177	▼0.0543	
Overall	▲0.0015	▼0.0023	▲0.0028	▼0.0043	▲0.0028	▼0.0043	

Table 5: Performance of Retrained Models with T14 and T15 Synthetic Samples

Class	Δ Balance	d Accuracy	$\Delta \mathbf{R}$	ecall	Δ F1-Score		
Class	T14	T15	T14	T15	T14	T15	
T1	▼0.0011	▼0.0004	▼0.0029	▼0.0017	▼0.0003	▲0.0005	
T2	▼0.0014	▼0.0039	▼0.0071	▼0.0097	▲0.0011	▼0.0030	
T3	▼0.0027	▼0.0027	▼0.0054	▼0.0043	▼0.0029	▼0.0063	
T4	▼0.0071	▲0.0069	▼0.0142	▲0.0142	▼0.0071	▲0.0037	
T5	▼0.0035	▼0.0044	▼0.0075	▼0.0094	▼0.0023	▼0.0027	
T6	▲0.0023	▼0.0027	▲0.0054	▼0.0054	▼0.0012	▼0.0042	
T7	▲0.0073	▲0.0023	▲0.0152	▲0.0051	▲0.0053	▼0.0007	
T8	▼0.0080	▼0.0104	▼0.0169	▼0.0212	▼0.0013	▼0.0101	
T9	▼0.0024	▼0.0031	▼0.0072	▼0.0072	▲0.0165	▲0.0047	
T10	▼0.0095	▲0.0016	▼0.0179	▲0.0030	▼0.0200	▲0.0039	
T11	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	
T12	▼0.0156	▲0.0001	▼0.0313	0.0000	▼0.0278	▲0.0122	
T13	0.0000	▼0.0111	0.0000	▼0.0222	0.0000	▼0.0243	
T14	▲0.0135	▼0.0259	▲0.0370	▼0.0519	▼0.1219	▼0.0388	
T15	▼0.0128	▲0.0219	▼0.0256	▲0.0513	▼0.0241	▼0.1079	
Overall	▼0.0026	▼0.0026	▼0.0049	▼0.0049	▼0.0049	▼0.0049	

These mixed results suggest that indiscriminate use of all generated synthetic samples may not consistently yield improvements. This observation led us to conclude that a selective approach to choosing text examples for synthetic data generation is necessary. Such selection consequently filters the generated samples to be appended for retraining the text classification model.

To evaluate the effectiveness of selection strategies, we first established a random selection baseline, followed by our proposed strategic selection methods. The following sections present these evaluations, beginning with the random selection baseline results.

B.1.3 SYNTHETIC SAMPLES PERFORMANCES WITH RANDOM EXAMPLES SELECTION

To establish a baseline for evaluating selection strategies, we implemented a ran-dom selection approach. For each target class, this baseline process randomly sam-ples a random number (between 1 and the size of the class pool) of examples with replacement, where selected examples subsequently go through the synthetic data generation process described at the beginning of Appendix B.1. We evaluated this baseline using 1 GPU hour fixed-time experiments to improve M0, maintaining consistency with the experimental settings described in Section 4.



Figure 15: Random Selection Results Distribution on Improving OBA (left) and CBA (Right).

The optimization process ran multiple iterations for 1 GPU hour, averaging 299 iterations per class, with 4487 sets of random selected examples tested in total. Figure 15 presents the performance distribution of models retrained using each ex-ample set. The overall balanced accuracy (OBA) improvement ranged from -0.9% to 0.2%, with an average of -0.3% and only 3.9% of retraining leading to positive OBA improvement. At the class level, the class balanced accuracy (CBA) improve-ment fluctuated between -5.8% and 4.7%, with an average of -2.2% and 9.1% of retraining achieving positive CBA improvement.

Table 6: Performance Comparison between Random Selection and Strategic Selections.

1116									
4447	Class	Ran	dom	Sliding	Window	Hierarch	ical SW	Genetic A	Algorithm
1117	Name	Overall	Class	Overall	Class	Overall	Class	Overall	Class
1118	T2	▲0.0012	▼0.0155	▲0.0030	▲0.0050	▲0.0034	▲0.0048	▲0.0009	▼0.0010
1119	T1	▲0.0009	▼0.0202	▲0.0028	▲0.0005	▲0.0029	▲0.0005	▲0.0018	▼0.0018
1120	T3	▲0.0014	▼0.0148	▲0.0030	▲0.0058	▲0.0027	▲0.0062	▲0.0029	▲0.0069
1120	T5	▲0.0002	▼0.0136	▲0.0032	▲0.0189	▲0.0034	▲0.0140	▲0.0010	▲0.0059
1121	T7	▲0.0018	▼0.0091	▲0.0036	▲0.0228	▲0.0034	▲0.0226	▲0.0012	▼0.0026
1122	T6	▲0.0016	▼0.0124	▲0.0035	▲0.0190	▲0.0030	▲0.0196	▲0.0015	▲0.0073
1123	T10	▼0.0006	▲0.0051	▲0.0034	▲0.0281	▲0.0044	▲0.0247	▲0.0027	▲0.0208
1104	T9	▲0.0016	▲0.0049	▲0.0036	▲0.0147	▲0.0027	▲0.0191	▲0.0026	▲0.0077
1124	T4	▲0.0018	▲0.0118	▲0.0029	▲0.0304	▲0.0033	▲0.0369	▲0.0036	▲0.0323
1125	T8	▲0.0012	▲0.0146	▲0.0030	▲0.0321	▲0.0029	▲0.0358	▲0.0020	▲0.0142
1126	T14	▲0.0009	▲0.0029	▲0.0023	▲0.0326	▲0.0029	▲0.0395	▲0.0019	▲0.0396
1197	T15	▲0.0005	▲0.0067	▲0.0033	▲0.0456	▲0.0034	▲0.0446	▲0.0037	▲0.0533
1121	T11	▲0.0020	▼0.0248	▲0.0030	▲0.0054	▲0.0030	▲0.0053	▲0.0039	▲0.0053
1128	T12	▲0.0014	▲0.0472	▲0.0037	▲0.0699	▲0.0032	▲0.0772	▲0.0036	▲0.0775
1129	T13	▲0.0006	▲0.0237	▲0.0030	▲0.0443	▲0.0037	▲0.0548	▲0.0034	▲0.0548

The best results achieved through random selection are presented in Table 6. While overall performance improvements were observed, the random selection struggled particularly with improving class performance for larger classes.

One set of plots showing the overall balanced accuracy (OBA) improvements over time for the 15 classes is presented in Figure 16, where the random selection benchmark results are shown in gray dash-lines. We can observe that random selection frequently achieves only marginal improvements over the original model, as evidenced in classes T7, T13, T14, and T15. Moreover, for T10, the random selection approach performed worse than the original model.



These findings motivated us to develop the three example selection strategies and conduct parameter studies for each, which we present in the following sections.







Figures 23a and 23b compare Level 0 Num Seg values of 8 and 2, when Window Size = Half. Num Seg = 8 shows more consistent good performance, leading to its selection.



- Population Size and Selection Size (PSSize)
 - Crossover Rate and Initial Mutation Rate (CMRate)

These parameters are represented by the four leftmost coordinates in the plot. We used the Hypervolume Indicator as the primary comparison metric, supplemented by the maximum values of both objectives in each of the 5 cross-validation folds.





1426 1427

1438

1439

1452 1453

1454

1455 1456

1457

1404

1405

1406

Figure 26: Genetic Algorithm Parameters Experiment Overview (PSSize: Population-Selection Size, CMRate: Crossover-Mutation Rate).

1428	Parallel Coordinates Plot for Genetic Algorithm Parameter Tests	Parallel Coordinates Plot for Genetic Algorithm Parameter Tests
1429	PS-Stack-Hear FO, HY E1, HY F2, HY F3, HY F4, HY MO (Maxi D Maxi D Maxi D Maxi D Maxi D Maxi T Maxi	PS-5tccCH-8tept D1 MF 11 MF 24 MF 31 MF 14 MF 04 NF 00450 D1450 24550 24
1430	0.3000 0.5000 0.0006 0.0009 0.0014 0.0006 0.8662 0.8552 0.8559 0.8558 0.8559 0.3250 0.3550 0.3500 0.3000 0.3000 0.3000 0.3000 0.3000 0.3000 0.3000 0.3000 0.3000 0.3000 0.3000 0.3000 0.3550 0.3550 0.8550 0.8550 0.3550 0.	1.5 0.9000 0.9000 0.0007 0.0006 0.0009 0.0014 0.0006 0.8550 0.8550 0.8558 0.3558 0.3250 0.8000 0.8000 0.8000 0.0009 0.0014 0.0006 0.8550 0.8550 0.8514 0.3250 0.8000 0.8000 0.8000 0.0009 0.0014 0.0006 0.8550 0.8514 0.005 0.3250 0.8514 0.8514 0.005 0
1431	0.7000 b.7000 b.0006 p.005 b.006 b.006 b.005 b.005 b.005 b.0651 b.0051 b	C.7000 b.7000 b.0006 0.0006 0.0006 0.0006 0.0007 0.0550 0.8556 0.8512 C.7000 b.7000 b.0005 0.0006 0.0016 0.0007 0.0550 0.8556 0.8512 D.0005 0.0006 0.0005 0.0006 0.0016 0.0007 0.0550 0.8556 0.8512 D.0005 0.0005 0.0006 0.0007 0.0550 0.8556 0.8512 D.0005 0.0005 0.0006 0.0007 0.0550 0.8556 0.8512 D.0005 0.0005 0.0005 0.0005 0.0007 0.0550 0.8556 0.8512 D.0005 0
1432	0.5000 0.5400 0.0005 0.0008 0.0013 0.8600 (8552 0.8550 0.5410 0.3350 0.3350 0.3350 0.3400 0.3550 0.3450 0.0007 2	2 0.5000 0.5000 0.0005 0.0006 0.0018 0000000000
1433	0.4000 4000 0.400 0.400 0.400 0.000 0.400 0.0550 0.8540 0.8540 0.8544 0.500 0.3850 0.0000 0.000 0.0000 0.000 0.0000 0.000 0.000 0.0000 0.000 0.000 0.0000 0.0000 0.000 0	LS 0.4000.4000 0.000 0.000 0.000 0.000 0.000 0.000 0.000 0.000 0.00000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.0000 0.000000
1434	0.2000 0.2000 0.2000 0.0007 0.0007 0.010 0.0007 0.011 0.000 0.8554 0.8549 0.8666 0.2950 0.3856 0.3550 0.3000 0.3856 0.3550 0.3000 0.3856 0.3550 0.3000 0.0006 0.0007 0.0007 0.0006 0.0007 0.000	0.2000 0.2000 0.2000 0.0003 0.0004 0.0007 0.001 0.0007 0.8594 0.8544 0.854 0.8542 0.8550 0.2950 0.3858 0.3550 0.30000 0.30000 0.30000 0.30000 0.3000 0.3000 0.3000 0.30000
1435		

1436(a) Performances of Population and Selection Sizes of
20 and 20.(b) Performances of Crossover and Initial Mutation
rates of 0.7 and 0.3.

Figure 27: Genetic Algorithm Parameters Experiments.

Figure 27a illustrates the performance when both population size and selection size
are set to 20. These results consistently outperform the alternative configuration of
population size 20 and selection size 10.

Figure 27b shows the performance with a crossover rate of 0.7 and an initial mutation rate of 0.3. This configuration demonstrates superior performance compared to the alternative rates of 0.9 and 0.1, respectively, when population and selection sizes are held constant.

Based on these experiments, we determined the optimal parameter set for the Genetic Algorithm strategy:

- Population Size = 20
- Selection Size = 20
- Crossover Rate = 0.7
 - Initial Mutation Rate = 0.3

1458 C ALGORITHMS

1460 C.1 EXAMPLES SUBSET SELECTION STRATEGIES 1461

Based on the parameter experiments reported in Appendix B, we finalized the workflows for the three example selection algorithms. To complement the description
provided in Section 3.3, we present here the detailed workflows for the Hierarchical Sliding Window (HSW) and Genetic Algorithm (GA) strategies.

Req	uire: $a_s - x$ and y range of each PCA plot; $n_s - list$ of number of segments for each level
	w_s – list of window sizes for each level; $data_syn$ – synthetic data; $class$ – the selected class
-	l – level of hierarchical sliding window
Ensu	ire: Best windows found on each level
1: 1	for each plot in PCA_plots do
2:	Initialize $l \leftarrow 0$
3:	Initialize selected_windows $\leftarrow \{a_s\}$
4:	while $l < \text{length of } n_{-s}$ and not terminated do
5:	$best_windows \leftarrow \{\}$
6:	for each area in selected_windows do
7:	Perform sliding window on area using $n_s[l]$ and $w_s[l]$
8:	for each window in sliding window do
9:	Retrieve data dots from window belonging to class
10:	$syn_samples \leftarrow AutoGeTS(data dots, data_syn)$
11:	Train classification model using syn_samples
12:	Evaluate model and compute performance metric $J'(W)$
13:	Record performance metric as the score for <i>window</i>
14:	end for
15:	Add window with maximum score to best_windows
16:	end for
17:	$selected_windows \leftarrow best_windows$
18:	$l \leftarrow l + 1$
19:	if termination condition met then
20:	Set terminated to True
21:	end if
22:	end while
23:	end for
~ ·	

1502 1503 1504

1505

C.2 ENSEMBLE MULTI-CLASS ALGORITHMS

This section presents detailed algorithms for the ensemble strategies outlined in Section 3.4. These algorithms are designed to address specific business requirements identified in Section 3.2.

to optimize the set of data examples used for synthetic data generation.

¹⁵¹⁰ Algorithm 3 addresses Requirement 1, focusing on improving the performance of underperforming classes. Algorithm 4 targets Requirement 2, aiming to enhance

Algor	ithm 2	GA Selec	tion-Gener	ation-Retraini	ng Approach						
Reaui	i re: n	– populati	on size: a _r		n generations: F	(S) - fit	ness score from objective				
fu	nction	; mutation	i - represent	ntation mutati	on function; <i>cro</i>	ssover –	representation crossover				
fu	function; p_{cro} – crossover probability; p_{mut} – mutation probability.										
Ensur	Ensure: The final individual(s) maximises the fitness score										
1: $P_0 \leftarrow$ randomly generated population of selected data examples of size n with priority value											
ba	ised ch	romosom(e representa	tion			through Arres CoTS and				
2: F	$0 \leftarrow \{$	$F(P_0[i])$	$i \in 1, \ldots$	$,n$ {Evaluat	e initial populatio	on intress	inrough AutoGe15 pro-				
3. G	:ss; (← 0 :	Generatio	n counter}								
4: while $G < g_{\text{max}}$ do											
5: Select individuals for the mating pool using Lexicase and Clustered tournament selection											
based on fitness score											
6: $P' \leftarrow$ Generate offspring using weight mapping <i>crossover</i> and adaptive polynomial <i>mutation</i>											
7.	With p	orobability	$p_{\rm cro}$ and $p_{\rm m}$	$\mathbf{F}(\mathbf{D}'[i]) \mid i \in \mathcal{I}$	1 ml (Evol	uoto now	nonulation				
7. 8∙	Comb	ate offsprei	and offspri	$I'(I [i]) \mid i \in$	$n_1, \ldots, n_f \in \mathbb{D}$ and $s: B \leftarrow P_C \sqcup P'$	uale new	population				
9:	Selec	t the next g	generation 1	P_{G+1} from R^{-1}	using elitism.						
10:	$G \leftarrow$	$G+1$ {In	crement ge	neration count	er}						
11: er	ıd whi	le	-		-						
12: R	eturn t	he best set	(s) of data e	examples in P	$_{G}$ based on fitnes	s score					
				Class Orde	rs by Requireme	זנ]				
1	5 11	T1 4/1					R1 Classes Order R2 Classes Order				
1	л т2	T1 4/1			T 5		R2 Classes Order R3 Classes Order				
1.	-										
1	3 71	T 8/3		T 5	T 6						
1	2 тз	T8/3		T 7	T 14	T 6					
1	1 т5	T5/2	T 14	T 14	T 15	T 5					
				•							
1	0 📑	T5/2	T1 T 13	T 13	T 12	T 13					
	о те	T 6	T 12 T 9	T3	T3 T12/1	T 8					
der											
s Ol	8 77	T7/4	T 11 T 2	T1	T1 T12/1	T 12					
Class											
0	7 710	T7/4	T 14 T 15/1	2 T 9/8 T 15	T4 T7	T 15	T 14				
	6 тв	Т9	T6 T1 5/1	2 T 9/8 T3	T10 T13	T 11	T 13				
	5 74	T 12	T 15 T5	T 15 T2 T	2 79 711	T3 T5/6	/9/14 15 11/5				
	4	115/13	9 110/1		3 11 14 15	9 15/0	/9/14 14 111/5				
	3 11	T1 5/13	T13 T1 0/1	T12 T6 T	11 T2 T2 T14	T2 T5/6	/9/14 T9 T10 T2				
				• • •							
	2 713	т 11 т2	T 5 T 11	T2 7 9 7	5 🕇 13 🕇 5 🕇 8	T4 T5/6	/9/14 T2 T12 T10				
		• ••• •••	T 2 T 2	-	c (77 (70 (70						
			2 13	4 5 7	0 1 1 8 1 9						
	Ŕ	82 317	312 313	314 315 316	·	in in	The the the				
		ب ۲	4 4 v	* * *	* * * * *	£ 65	43 43 43 4 <u>3</u>				
				Re	equirement						



overall classification performance. Algorithm 5 addresses Requirement 3, which
 prioritizes improving the performance of a specific important class.

1569 Each algorithm iteratively utilizes the AutoGeTS process, incorporating insights 1570 from our experimental results, including the strategy-objective combinations from 1571 Table 3 and class relationships and class order from Figures 6b and 6a. In Figure 6a, 1572 each row represents the application of AutoGeTS to one of the fifteen classes, while 1573 each column shows the corresponding improvement in class balanced accuracy for 1574 that particular class. The resulting class orderings for the three requirements are il-1575 lustrated in Figure 28. The blue column depicts the class ordering for Requirement 1576 1577 1 (improving all individual class performance), where the bottom-most dot repre-1578 sents the first-ordered class and the top-most dot indicates the last-ordered class. 1579 The green column shows the class ordering for Requirement 2 (improving overall 1580 performance), while the orange columns display the class ordering for Requirement 1581 3 (improving a specific important class), with each R3 column corresponding to a 1582 designated important class. 1583

Alg	orithm 3 Requirement 1: Improve Bad Performing Classes
Rec	uire: Classes with performance metrics, AutoGeTS process, GA parameters
Ens	ure: Improved classification model for bad-performing classes
1:	Sort classes by class size in ascending order
2:	for Iteration <i>i</i> , each class C_i with class balanced accuracy < 0.8 do
3:	Select strategy-objective that best improves C_i balanced accuracy from Lookup Table 3
4:	Apply AutoGeTS to C_i
5:	if improvement achieved then
6:	Record maximum improvement and corresponding model m for C_i
7:	if exist M maintain improvement in each $C_{n < i}$ balanced accuracy by at least 50% then
8:	Select model m from M that best improves C_i 's balanced accuracy
9:	else
10:	Terminate algorithm
11:	end if
12:	Append the synthetic sample from m to the training set
13:	eise Terminete elgorithm
14:	and if
15. 16·	end for
Alg	orithm 4 Requirement 2: Improve Overall Performance
Rea	uire: Classes with performance metrics, AutoGeTS process, GA parameters
Ens	sure: Improved overall classification performance
1:	while overall performance can be improved do
2:	Select a class C in descending order of improving overall performance based on figure 6b.
3:	Select strategy-objective combination that best improves global metrics for C from Looku
	Table 3
4:	Apply AutoGeTS to C
5:	Record maximum improvement and corresponding model m for C
6:	Append the synthetic sample from m to the training set
7:	if overall performance not improved then
8:	Terminate algorithm
9:	ena II
1()	ena while

Algo	rithm 5 Requirement 3: Improve Important Class
Requ	ire: Classes with performances, AutoGeTS process, GA parameters, Important class IC
Ensu	re: Improved important class performance
1: Io	dentify related classes RC of IC from figure 6a.
2: S	ort IC and RC in descending order of IC improvement according to figure 6a.
3: fo	or iteration <i>i</i> , each class IC_i or RC in the order do
4:	Apply AutoGeTS to IC_i or RC with the strategy-objective combination that best improves
5.	Record maximum improvement in IC's class performance and corresponding model m
5: 6:	Append the synthetic sample from m to the training set
0. 7·	if <i>IC</i> performance not improved then
8:	Terminate algorithm
9:	end if
10: e	nd for
D	Fixed-Time Experiments
D.1	FIXED-TIME EXPERIMENTS: IMPROVING LOCAL VS GLOBAL METRIC
Follo	owing our analysis of the Performance Improvement Overview in Section 4.2
and	before comparing strategies and objectives in Section 4.3 we investigated
what	har AutoCoTS could simultaneously improve both class specific and overall
witer	inci Autoberts could simulateously improve bour class-specific and overall
perfo	ormance, and to what extent these goals might be contradictory. To this end,
we c	ompared models trained with synthetic data that achieved maximum improve-
ment	ts in either the local metric (Class Balanced Accuracy) or the global metric
(Ove	erall Balanced Accuracy) for all 15 classes.
` 	
Anal	lysis of Figure 29 reveals significant insights into AutoGeTS's performance.
Whe	n optimizing for local metrics, 11 out of 15 classes (excluding T1, T7, T8, and
T14)	showed improvements without negatively impacting the global metric. Only
T7 a	and T8 exhibited relatively larger decreases (-0.1%) in global performance
whei	n local performance was maximized. Similarly when optimizing for global
	in octa performance was maximized. Similarly, when optimizing for global
metr	ics, 11 out of 15 classes (excluding 11, 19, 111, and 114) demonstrated im-
prov	ements without compromising local performance. In this case, only 19 showed
a rel	atively larger decrease (-1%) in local performance when global performance
was	maximized. These observations demonstrate AutoGeTS's capability to simul-
taneo	ously improve both local and global metrics in the majority of cases. confirming
its of	fectiveness in addressing both class-specific and overall performance goals
113 01	need veness in addressing bour class-specific and overall performance goals.
How	ever, the instances where trade-offs occurred between local and global perfor-
man	ce suggest the need for future research into advanced optimization methods.
Such	research could explore both example selection strategies and objective func-
tions	anable of ontimizing multiple objectives simultaneously notantially alimi
uons	capable of optimizing multiple objectives simultaneously, potentially elimi-
natir	ng these trade-offs and further enhancing AutoGeTS's performance across all
class	es.
D.2	FIXED-TIME EXPERIMENTS: PERFORMANCE TRAJECTORIES OVER RETRAINING TIME
Follo	owing the comparison of the three example selection strategies in Section 4.3.
we f	further analyzed their performance improvements with respect to retraining



Under review as a conference paper at ICLR 2025



time. Four experiments were conducted, each maximizing a different metric: Class
Recall, Class Balanced Accuracy, Overall Balanced Accuracy, and Overall F1Score. The three example search strategies and 15 classes served as independent
variables, with a constraint of 1 hour total GPU running time.

Figure 30 compares improvements in Class Recall when maximizing Class Recall is the optimization objective. Figure 31 compares improvements in Class Balanced Accuracy when maximizing Class Balanced Accuracy is the optimization objective. Figure 32 compares improvements in Overall Balanced Accuracy when maximiz-ing Overall Balanced Accuracy is the optimization objective. Figure 33 compares improvements in Overall F1-Score when maximizing Overall F1-Score is the opti-mization objective.

For these class-specific metrics, we observed that HSW often achieves its best or near-best improvements within the first 1/3 of training time, especially for classes where it performs best. GA typically reaches its peak performance within the first 1/3 of retraining time, except for T14 and T15. However, for these classes, GA outperforms other strategies even before reaching its best performance. These observations confirm the computational efficiency of HSW and GA strategies.

SW also often achieves its best performance around 1/3 of retraining time. But for
classes where it outperforms other strategies (e.g., T5 and T10), SW only surpasses
others quite late, starting from around 2/3 or even 4/5 of retraining time. This
reflects both its advantage in avoiding plateaus and its computational inefficiency
due to its brute-force nature.

Figure 33 compares improvements in Overall F1-Score when maximizing OverallF1-Score is the optimization objective.





Figure 30: Fixed Time, Class Recall as Objective, Comparing on Class Recall Improvement: each time series plot has six lines. SW max, HSW max, GA max, SW avg, HSW avg, GA avg.



Figure 31: Fixed Time, Class Balanced Accuracy as Objective, Comparing on Class Balanced Accuracy Improvement: each time series plot has six lines. SW max, HSW max, GA max, SW avg, HSW avg, GA avg.



Figure 32: Fixed Time, Overall Balanced Accuracy as Objective, Comparing on Overall Balanced
Accuracy Improvement: each time series plot has six lines. SW max, HSW max, GA max, SW avg,
HSW avg, GA avg.



Figure 33: Fixed Time, Overall F1-Score as Objective, Comparing on Overall F1-Score Improvement: each time series plot has six lines. SW max, HSW max, GA max, SW avg, HSW avg, GA avg.

For these global metrics, we observed that all three strategies often required more than 1/3 of retraining time to reach their best or near-best performance, especially for Overall Balanced Accuracy. This suggests improving global metrics would be more computationally complex than improving local metrics for all three strategies, as it would require a more synergistic effect.

When improving global metrics, GA often leads in performance gain from early on (around 1/3 retraining time) for classes where it obtains the largest improvement. SW and HSW show more variable timing in achieving leading performance, ranging from early to late in the retraining process. GA's performance pattern con-

firms its evolutionary nature, demonstrating an ability to maintain and exploit good solutions once found.

BENCHMARK DATA AUGMENTATION COMPARISON EXPERIMENTS Ε

To provide an additional baseline for evaluating the AutoGeTS workflow, we im-plemented traditional data augmentation methods, alongside the LLM-based ap-proach, as the generator in the synthetic data generation process. Specifically, we utilized the Easy Data Augmentation (EDA) tool (Wei & Zou (2019)), which incor-porates techniques such as synonym replacement, random insertion, random swap, and random deletion. The evaluation followed our established experimental pro-tocol of 1 GPU hour fixed-time experiments to improve M0, employing Random, SW, HSW, and GA selection strategies, each executed separately with either max-imizing Overall Balanced Accuracy (OBA) or Class Balanced Accuracy (CBA) as the optimization objective. All experimental parameters remained consistent with those described in Section 4.

E.1 BEST PERFORMANCE COMPARISON: OBA & CBA IMPROVEMENTS

Table 7: Overall Balanced Accuracy (Global) Comparison between EDA and GPT-3.5 as Generator.

2022		D	1		TT 7' 1	TT' 1	. 1000		1 11
2022	Class	Ran	dom	Sliding	Window	Hierarch	iical SW	Genetic A	Algorithm
2023	Name	EDA	GPT-3.5	EDA	GPT-3.5	EDA	GPT-3.5	EDA	GPT-3.5
2024	T2	▲0.0001	▲0.0012	▲0.0012	▲0.0030	▲0.0014	▲0.0034	▲0.0004	▲0.0009
	T1	▲0.0006	▲0.0009	▲0.0017	▲0.0028	▲0.0024	▲0.0029	▲0.0012	▲0.0018
2025	T3	▲0.0006	▲0.0014	▲0.0011	▲0.0030	▲0.0018	▲0.0027	▲0.0010	▲0.0029
2026	T5	▲0.0003	▲0.0002	▲0.0014	▲0.0032	▲0.0011	▲0.0034	▲0.0006	▲0.0010
2027	T7	▲0.0006	▲0.0018	▲0.0009	▲0.0036	▲0.0015	▲0.0034	▲0.0006	▲0.0012
2028	T6	▲0.0010	▲0.0016	▲0.0010	▲0.0035	▲0.0013	▲0.0030	▲0.0013	▲0.0015
2020	T10	▼0.0011	▼0.0006	▲0.0008	▲0.0034	▲0.0015	▲0.0044	▲0.0016	▲0.0027
2029	T9	▲0.0001	▲0.0016	▲0.0019	▲0.0036	▲0.0014	▲0.0027	▲0.0011	▲0.0026
2030	T4	▲0.0010	▲0.0018	▲0.0005	▲0.0029	▲0.0019	▲0.0033	▲0.0018	▲0.0036
2031	T8	▲0.0006	▲0.0012	▲0.0008	▲0.0030	▲0.0011	▲0.0029	▲0.0014	▲0.0020
2022	T14	▲0.0006	▲0.0009	▲0.0014	▲0.0023	▲0.0019	▲0.0029	▲0.0003	▲0.0019
2032	T15	▲0.0001	▲0.0005	▲0.0011	▲0.0033	▲0.0013	▲0.0034	▲0.0022	▲0.0037
2033	T11	▲0.0012	▲0.0020	▲0.0014	▲0.0030	▲0.0014	▲0.0030	▲0.0022	▲0.0039
2034	T12	▲0.0010	▲0.0014	▲0.0011	▲0.0037	▲0.0021	▲0.0032	▲0.0020	▲0.0036
2035	T13	▲0.0004	▲0.0006	▲0.0013	▲0.0030	▲0.0011	▲0.0037	▲0.0022	▲0.0034

Table 8: Class Balanced Accuracy (Local) Comparison between EDA and GPT-3.5 as Generator.

2000									
2039	Class	Random		Sliding Window		Hierarchical SW		Genetic Algorithm	
2040	Name	EDA	GPT-3.5	EDA	GPT-3.5	EDA	GPT-3.5	EDA	GPT-3.5
2040	T2	▼0.0212	▼0.0155	▲0.0003	▲0.0050	▲0.0003	▲0.0048	▼0.0011	▼0.0010
2041	T1	▼0.0219	▼0.0202	▼0.0030	▲0.0005	▼0.0027	▲0.0005	▼0.0028	▼0.0018
2042	T3	▼0.0215	▼0.0148	▲0.0015	▲0.0058	▲0.0007	▲0.0062	▼0.0012	▲0.0069
2043	T5	▲0.0001	▼0.0136	▲0.0139	▲0.0189	▲0.0149	▲0.0140	▲0.0048	▲0.0059
2044	T7	▼0.0082	▼0.0091	▲0.0191	▲0.0228	▲0.0180	▲0.0226	▲0.0056	▼0.0026
2044	T6	▼0.0055	▼0.0124	▲0.0122	▲0.0190	▲0.0195	▲0.0196	▲0.0093	▲0.0073
2045	T10	▲0.0013	▲0.0051	▲0.0244	▲0.0281	▲0.0255	▲0.0247	▲0.0203	▲0.0208
2046	T9	▼0.0174	▲0.0049	▲0.0041	▲0.0147	▲0.0071	▲0.0191	▼0.0023	▲0.0077
2047	T4	▲0.0138	▲0.0118	▲0.0293	▲0.0304	▲0.0272	▲0.0369	▲0.0340	▲0.0323
2041	T8	▼0.0000	▲0.0146	▲0.0181	▲0.0321	▲0.0285	▲0.0358	▲0.0242	▲0.0142
2048	T14	▼0.0014	▲0.0029	▲0.0153	▲0.0326	▲0.0221	▲0.0395	▲0.0185	▲0.0396
2049	T15	▲0.0038	▲0.0067	▲0.0245	▲0.0456	▲0.0409	▲0.0446	▲0.0199	▲0.0533
2050	T11	▼0.0143	▼0.0248	▲0.0051	▲0.0054	▲0.0052	▲0.0053	▲0.0051	▲0.0053
2051	T12	▲0.0405	▲0.0472	▲0.0656	▲0.0699	▲0.0705	▲0.0772	▲0.0709	▲0.0775
2001	T13	▲0.0204	▲0.0237	▲0.0415	▲0.0443	▲0.0514	▲0.0548	▲0.0392	▲0.0548

As our deployment strategy selects the best-performing model for each target class based on specified objectives, we analyze the maximum OBA and CBA improvements achieved and compare them with the AutoGeTS performance presented in Section 4.2. Tables 7 and 8 present these comparative results.

2057 The comparison between EDA and LLM (GPT-3.5) approaches reveals consis-2058 tent performance advantages for the LLM-based AutoGeTS workflow across all 2059 selection strategies. With the random selection, EDA's best-performing models 2060 performed below LLM's by margins of 0.06% in OBA and 0.25% in CBA, un-2061 2062 derperforming in 15 and 11 classes respectively. This performance gap widened 2063 with strategic selection methods: under SW, EDA performed 0.20% lower in OBA 2064 and 0.69% lower in CBA than LLM, with inferior results across all 15 classes for 2065 both metrics. Similarly, with HSW, EDA showed 0.17% lower OBA and 0.51% 2066 lower CBA on average, underperforming in 15 and 13 classes respectively. The GA 2067 strategy yielded comparable results, with EDA's best-performing models averaging 2068 2069 0.11% below LLM in OBA and 0.51% below in CBA, showing lower performance 2070 in 15 and 11 classes respectively.

2071 2072

E.2 COMPARISON OF OBA IMPROVEMENTS OVER TIME

We extend our analysis to examine the temporal progression of overall balanced accuracy (OBA) improvements across all 15 classes, comparing EDA and LLM (GPT-3.5) results within the 1-GPU hour experimental constraint.

Figure 34 presents these comparisons, displaying the results of four selection strategies (SW, HSW, GA, and random selection) for both approaches. The EDA results are depicted with solid lines, while the corresponding LLM results are shown with dotted or dashed lines of the same color. Across all 15 classes, each EDA trajectory (solid line) mostly falls below its LLM counterpart (dotted or dashed line), demonstrating the superior overall performance of the LLM-based AutoGeTS workflow.



Figure 34: Fixed 1 Hour GPU time (x-axis, in seconds), Comparing on OBA Improvement (y-axis): solid lines are from EDA and dotted lines are from GPT-3.5.