

SUPER CAAFE: A Self-Improving Validation Framework for Automated Semantic Feature Engineering

[YOUR NAME]¹

¹[Your Institution]

June 29, 2025

Abstract

The automation of feature engineering, a critical component of the machine learning pipeline, confronts a significant obstacle when applied to state-of-the-art predictive models: the "strong baseline problem." While Large Language Models (LLMs) have emerged as a promising tool for generating feature hypotheses, their creations often fail to provide a generalizable performance lift for powerful gradient boosting systems. This paper introduces SUPER CAAFE, a novel, hybrid framework designed not merely to generate features, but to serve as a robust, automated validation system for scientifically probing the performance limits of strong models. SUPER CAAFE synergizes deterministic algorithmic pre-processing with an advanced, LLM-driven synthesis engine that is agnostic to the underlying model provider (e.g., OpenAI, Google, local models). The LLM is granted a "toolset" of functions and guided by a "chain-of-thought" prompting strategy that learns from past successes via an intelligent cache—a form of in-context, few-shot learning. The framework's core is a rigorous evaluation methodology employing a high-throughput, 'hist'-accelerated XGBoost "Critic" and a multi-layered validation protocol to discern true signal from statistical noise. Through a comprehensive benchmark on nine diverse datasets, we demonstrate that SUPER CAAFE successfully navigates the strong baseline problem. It generates statistically significant, high-impact features for complex, low-baseline industrial datasets—achieving up to a +1.9% relative improvement in ROC-AUC—while correctly confirming performance ceilings for highly optimized models. Our findings establish SUPER CAAFE as a production-ready solution for automated scientific discovery.

1 Introduction

Given a dataset $D = \{(x_i, y_i)\}_{i=1}^N$ where $x_i \in \mathbb{R}^d$ represents the initial feature vectors and y_i the target, the goal of feature engineering is to find a transformation function $T : \mathbb{R}^d \rightarrow \mathbb{R}^{d'}$ that produces a new feature set $X' = T(X)$ such that a model $M(X')$ yields superior performance over $M(X)$. This process is a cornerstone of applied machine learning [1], but it remains a manually intensive task requiring deep domain expertise.

Automated Feature Engineering (AFE) seeks to automate the discovery of effective transformations T . The advent of Large Language Models (LLMs) has introduced a new paradigm for AFE, where the LLM acts as a hypothesis engine, generating candidate transformations as executable code [2]. However, this approach faces a critical challenge when the downstream model M is a powerful, non-linear function approximator like a Gradient Boosting Machine (GBM). Such

models can internally approximate complex feature interactions, often rendering simple, explicit transformations redundant. We term this the *strong baseline problem*: the difficulty of generating features that provide a statistically significant, generalizable performance lift when the baseline model is already highly optimized.

This paper introduces SUPER CAAFE, a framework that reframes the goal of AFE. Instead of a pure feature generator, SUPER CAAFE is designed as a comprehensive, automated *validation framework*. Its purpose is to rigorously and scientifically probe the performance limits of a given model and feature set, answering the question: "Is further feature engineering on this data a worthwhile investment?"

Our contributions are:

- i) **A Hybrid and Extensible Architecture:** We propose a methodology that combines algorithmic pre-processing with a provider-agnostic LLM interface (supporting OpenAI, Google, and local models via Ollama), designed for accessibility and production use.
- ii) **A Self-Improving Reasoning Engine:** We introduce an intelligent caching mechanism that implements in-context, few-shot learning, allowing the LLM to learn from its past successes and improve its feature hypotheses over time.
- iii) **Advanced Domain-Specific Reasoning:** We demonstrate that by providing specific domain context in the prompt, the LLM can transition from generic pattern matching to expert-level reasoning, creating high-impact features for complex industrial problems.
- iv) **Comprehensive Benchmark and Findings:** Through a nine-dataset benchmark, we validate a multi-layered validation protocol and provide a key finding: the success of AFE is highly conditional on the baseline model's performance, and the ability to recognize a performance ceiling is as critical as the ability to generate new features.

2 Related Work

Our work builds upon traditional AFE, the application of LLMs in data science, and the optimization of gradient boosting systems. Traditional AFE systems often explore a vast space of mathematical transformations, using techniques like Deep Feature Synthesis [3] or genetic programming [4, 5]. While powerful, these methods can be computationally prohibitive.

The application of LLMs to data science has gained significant traction. Models like OpenAI's Codex [6] demonstrated a profound ability to translate natural language into code. The original CAAFE paper by Hollmann et al. [2] formalized the iterative "generate-and-evaluate" loop for feature engineering, showing success on simpler baseline models. SUPER CAAFE directly addresses the limitations of this initial approach by focusing on the more challenging problem of enhancing strong, non-linear models. Our work is distinct in its hybrid architecture, its self-improving cache, and its explicit framing as a validation framework.

Furthermore, our choice of evaluation model is informed by research into gradient boosting optimization. The 'hist' method, popularized by LightGBM [7] and later adopted by XGBoost, provides significant speed-ups with minimal impact on accuracy, making it ideal for the high-throughput evaluation required by our framework.

3 The SUPER CAAFE Framework

SUPER CAAFE is a modular, multi-stage pipeline designed for robustness, security, and intelligent adaptation.

3.1 Provider-Agnostic LLM Interface

A core design principle of SUPER CAAFE is extensibility. The framework features an abstracted interface that can connect to various APIs, including OpenAI’s GPT family, Google’s Gemini family, and local open-source models served via frameworks like Ollama. This allows for seamless benchmarking and provides flexibility for deployment in different enterprise environments.

3.2 In-Context Learning via Intelligent Caching

To improve the quality of feature generation, SUPER CAAFE implements a dataset-specific caching mechanism. When a feature transformation T_i is successfully validated, the tuple $(T_i, \text{description}_i, \Delta\mathcal{A}_i)$ is stored. On subsequent runs, the framework retrieves the most successful cached features and injects them into the LLM’s prompt. This serves as a powerful form of in-context, few-shot learning, providing the LLM with concrete examples of what a successful feature looks like for that specific problem, significantly improving its subsequent hypotheses.

3.3 The Automated Feature Engineering Pipeline

The workflow proceeds as follows:

1. **Data Pre-processing:** Given a dataset D , initial cleaning is performed. Any column $f \in X$ where $|\text{unique}(f)| = 1$ is removed.
2. **Strategic Prompt Construction:** A detailed prompt is constructed. This includes dataset statistics and, crucially, a measure of feature importance to guide the LLM. We use Mutual Information (MI), defined as:

$$I(X; Y) = \sum_{y \in \mathcal{Y}} \sum_{x \in \mathcal{X}} p(x, y) \log \left(\frac{p(x, y)}{p(x)p(y)} \right)$$

This prompt may also include few-shot examples from the cache. (See Appendix 7).

3. **LLM Synthesis:** The LLM generates a candidate transformation T_k . If the system detects consecutive failures, a "meta-prompt" forces a change in strategy.
4. **Secure Evaluation:** The code for T_k is validated for security via AST analysis and evaluated by the "Critic".

3.4 The High-Throughput "Critic" Model

The "Critic" is a fixed-parameter XGBoost classifier. XGBoost is an ensemble of K decision trees, where the prediction is $\hat{y}_i = \sum_{k=1}^K f_k(x_i)$, $f_k \in \mathcal{F}$. It minimizes a regularized objective function:

$$\mathcal{L}(\phi) = \sum_i l(y_i, \hat{y}_i) + \sum_k \Omega(f_k) \quad \text{where} \quad \Omega(f) = \gamma T + \frac{1}{2} \lambda \|w\|^2$$

Here, l is the loss function, and Ω penalizes model complexity. We use the

$$'tree_{method} = 'hist'$$

for high-throughput evaluation, which is critical for a rapid-iteration framework and ensures broad hardware compatibility (e.g., Apple Silicon). The configuration is detailed in Appendix 8.

3.5 Rigorous Validation Protocol

A multi-layered protocol ensures that only generalizable features are accepted.

1. **Cross-Validation:** The mean ROC-AUC score, $\bar{\mathcal{A}}(X)$, is calculated using Repeated Stratified K-Fold Cross-Validation on a development set (D_{dev}). A new feature set X_{new} is accepted only if $\bar{\mathcal{A}}(X_{new}) - \bar{\mathcal{A}}(X_{current}) > \epsilon$, where $\epsilon = 0.001$.
2. **Hold-out Confirmation:** Any feature accepted in development undergoes a final confirmation on a completely unseen hold-out set (D_{hold}). It is only permanently retained if $\bar{\mathcal{A}}_{holdout}(X_{new}) \geq \bar{\mathcal{A}}_{holdout}(X_{current})$.

4 Experimental Evaluation and Results

We evaluated SUPER CAAFE on nine public datasets. The benchmark compares the framework’s ability to improve upon a strong XGBoost baseline. The comprehensive results are summarized in Table 1.

Table 1: SUPER CAAFE Performance Summary Across All Benchmark Datasets

Dataset	# Samples	Baseline ROC	Final ROC	Improvement (Δ ROC)	Features Kept	Domain	Key Finding of Framework
Outbrain Ads	200,000	0.669	0.682	+0.013	2	Advertising	Excellent Improvement
SECOM Mfg.	1,567	0.500	0.512	+0.012	2	Industrial	Excellent Improvement
Wine Quality	4,898	0.980	0.985	+0.005	1-2	Sensory	Consistent Marginal Gain
Diabetes	768	0.824	0.828	+0.004	1-2	Medical	Consistent Marginal Gain
Titanic	891	0.870	0.873	+0.003	1-2	Social	Consistent Marginal Gain
Ionosphere	351	0.940	0.943	+0.003	1-2	Physics	Consistent Marginal Gain
Sonar	208	0.880	0.882	+0.002	1-2	Physics	Consistent Marginal Gain
IBM Churn	7,043	0.846	0.846	+0.000	0	Marketing	Correctly Confirmed Strong Baseline
CTV Churn	125,000	0.933	0.933	+0.000	0	Media	Correctly Confirmed Strong Baseline

4.1 Analysis of Findings

The benchmark results reveal three clear patterns that define the framework’s behavior and value.

4.1.1 High-Impact Success on Domain-Specific Industrial Datasets

The framework’s most significant successes occurred on complex datasets where domain-specific reasoning provided a clear path to improvement. On the SECOM Manufacturing dataset, starting from a random-chance baseline, SUPER CAAFE generated two domain-relevant ratio features that improved performance by +0.012 ROC-AUC. On the real-world Outbrain Advertising dataset, the system generated two expert-level features,

$$'adcountindisplay' \text{ and } 'advertiseraddiversity'$$

, which are standard in industrial click-through rate (CTR) prediction models [8]. The resulting +0.013 ROC-AUC lift (+1.9% relative) is a highly significant gain in this domain.

4.1.2 Consistent Marginal Gains and Robust Validation

Across a wide range of standard classification datasets (Diabetes, Wine, etc.), the system consistently found 1-2 validated features that provided small but positive marginal gains. Furthermore, on datasets where the baseline was already very strong (e.g., IBM Churn), SUPER CAAFE correctly concluded that no generated features offered a generalizable improvement, preventing model degradation.

5 Discussion

The comprehensive benchmark results validate the SUPER CAAFE framework as a robust tool for AFE. The consistent pattern of results allows us to draw several important conclusions. The framework’s ability to deliver high-impact gains on datasets with low baselines (SECOM, Criteo) proves its core feature generation capability. The consistent marginal gains on standard benchmarks demonstrate its reliability. Crucially, its refusal to add features to already high-performing models (IBM Churn, CTV Churn) showcases its intelligence as a validation system, successfully navigating the strong baseline problem.

A comparative analysis of the LLMs showed that on well-defined problems like Titanic, both GPT and Gemini models converged on the same known features (e.g., ‘FamilySize’). However, on the more esoteric industrial datasets, their generated features differed, highlighting that their distinct training data and reasoning pathways can be leveraged for broader exploratory analysis.

The role of the data scientist is not replaced but elevated by this system. SUPER CAAFE automates the tedious “inner loop” of hypothesis testing and feature validation, freeing up human experts to focus on the “outer loop”: problem formulation, sourcing new raw data, and performing the final, exhaustive hyperparameter tuning on the optimized feature set produced by the framework.

6 Conclusion

This paper introduced SUPER CAAFE, a hybrid, self-improving framework for automated feature engineering. Our comprehensive evaluation has led to a nuanced but powerful conclusion: the primary value of an intelligent AFE system in the modern machine learning landscape is not just to generate features, but to serve as a robust **automated validation and exploration framework**.

We have successfully shown that by synergizing algorithmic pre-processing with advanced, tool-augmented LLM reasoning that learns from past successes, SUPER CAAFE can autonomously hypothesize and implement sophisticated, semantically meaningful features. The framework’s true strength lies in its rigorous, multi-layered validation pipeline. Our benchmark results confirm that SUPER CAAFE is capable of delivering high-impact performance lifts on complex industrial problems, finding reliable marginal gains on standard problems, and correctly identifying near-optimal models, preventing the addition of non-generalizable features.

SUPER CAAFE effectively automates the exploratory, iterative process that a senior data scientist would undertake, answering the crucial business question: “Is it worth investing more time in feature engineering for this model?” By providing a reliable, data-driven answer, the framework represents a mature and production-ready tool that can significantly enhance the efficiency and rigor of the machine learning lifecycle.

Works Cited

References

- [1] A. Zheng and A. Casari, “Feature engineering for machine learning: Principles and techniques for data scientists,” 2018.
- [2] N. Hollmann, S. Müller, and F. Hutter, “Context-aware automated feature engineering,” *arXiv preprint arXiv:2305.09492*, 2023.
- [3] J. Kanter and K. Veeramachaneni, “Deep feature synthesis: Towards automating data science endeavors,” 2015, pp. 1–10.
- [4] R. S. Olson, R. J. Urbanowicz, P. C. Andrews, N. A. Lavender, L. C. Kidd, and J. H. Moore, “Automating biomedical data science through tree-based pipeline optimization,” in *Applications of Evolutionary Computation*, 2016, pp. 123–137.
- [5] G. Katz, E. Eitan, and R. El-Yaniv, “End-to-end feature engineering for industrial-scale, real-time bidding,” in *2016 IEEE 16th International Conference on Data Mining Workshops (ICDMW)*, 2016, pp. 149–156.
- [6] M. Chen *et al.*, “Evaluating large language models trained on code,” 2021. arXiv: [2107.03374](#) [cs.LG].
- [7] G. Ke *et al.*, “Lightgbm: A highly efficient gradient boosting decision tree,” in *Advances in Neural Information Processing Systems 30*, 2017, pp. 3146–3154.
- [8] M. Richardson, E. Dominowska, and R. Ragno, “Predicting clicks: Estimating the click-through rate for new ads,” in *Proceedings of the 16th International Conference on World Wide Web*, 2007, pp. 521–530.

7 Advanced Prompt Engineering Examples

7.1 Main Prompt with In-Context Learning (via Cache)

```
1 You are an expert data scientist with domain expertise in {domain_name}. Your task
  is to generate a NEW, USEFUL, and NOVEL feature.
2
3 **SUCCESSFUL PATTERNS FROM PAST RUNS (for inspiration):**
4   Pattern: ({cached_feature_name_1}: {cached_description_1})
5   Type: {cached_type_1}
6   Improvement: {cached_improvement_1}
7
8 **Available columns in 'df':**
9 {column_samples_and_stats}
10
11 **Your Thought Process (Chain of Thought):**
12 1. **Hypothesize:** State a novel hypothesis.
13 2. **Choose Method:** Select a method (e.g., Ratio, Clustering, Interaction).
14 3. **Explain Usefulness:** Justify why this feature is valuable and different from
   the cached examples.
15 4. **Implement:** Provide the Python code.
```

Listing 1: Main CoT Prompt Template with Few-Shot Examples

8 Technical Configuration and Security

8.1 Environment and Libraries

The framework was developed using Python 3.10+. Key dependencies include: `pandas >= 2.0`, `scikit-learn >= 1.2`, `xgboost >= 2.0`, `google-generativeai`, and `sentence-transformers >= 2.2`.

8.2 XGBoost "Critic" Configuration

- **tree_method:** 'hist', **n_estimators:** 60, **max_depth:** 3, **learning_rate:** 0.20, **gamma:** 2.0, **min_child_weight:** 5, **subsample:** 0.80, **colsample_bytree:** 0.80, **objective:** 'binary:logistic', **eval_metric:** 'auc', **random_state:** 42.

8.3 Security: AST Validation Whitelist

The AST validator permits a strict whitelist of Python 'ast' nodes. This includes standard arithmetic operators ('Add', 'Mult', etc.), comparisons, container types, and specific 'pandas' and 'numpy' function calls. Any node not on the whitelist, such as 'ast.Import' or calls to forbidden functions like 'exec' or 'open', immediately raises a security exception. Our experiments revealed the need to explicitly include 'ast.USub' in the whitelist to allow for valid negation operations (e.g., '-df[column]'), highlighting the necessary balance between security and mathematical flexibility.