

# PROMPTFE: AUTOMATED FEATURE ENGINEERING BY PROMPTING

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

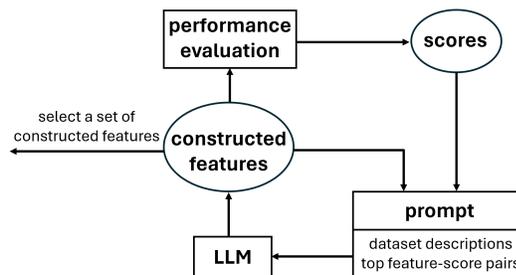
Automated feature engineering (AutoFE) liberates data scientists from the burden of manual feature construction. The semantic information of datasets contains rich context information for feature engineering but has been underutilized in many existing AutoFE works. We present PromptFE, a novel AutoFE framework that leverages large language models (LLMs) to automatically construct features in a compact string format and generate semantic explanations based on dataset descriptions. By learning the performance of constructed features in context, the LLM iteratively improves feature construction. We demonstrate through experiments on real-world datasets the superior performance of PromptFE over state-of-the-art AutoFE methods. We verify the impact of dataset semantic information and provide comprehensive study on the LLM-based feature construction process.

## 1 INTRODUCTION

Tabular data, a form of structured data comprising instances and attributes, have extensive use across a broad range of domains including credit assessment, market prediction, and quality control. Traditional machine learning models, especially tree-based models (Breiman, 2001; Ke et al., 2017), have strong performance on tabular datasets of small and medium sizes (Grinsztajn et al., 2022) and good interpretability. Feature engineering refers to the construction of features to enhance the performance of downstream models, which is crucial for traditional ML models as new features extract useful information for target prediction by capturing complex non-linear relationships. Feature engineering by hand demands domain expertise to relieve significant human labor.

Automated feature engineering (AutoFE) employs meta algorithms and models to automate feature engineering process for performance comparable to domain experts. Prior approaches like (Zhu et al., 2022a;b; Zhang et al., 2023) construct and evaluate enormous features in a trial-and-error manner. While some methods learn to optimize the utility of features during the FE process, they do not utilize domain knowledge to guide feature search. The need to search features from scratch for new datasets and downstream models hampers their efficacy and efficiency. Furthermore, these methods cannot offer explanation of the engineered features, undermining the interpretability.

The text descriptions of tabular datasets provide rich context for feature engineering. Domain experts consult attribute descriptions to select relevant feature attributes and compute new features useful for target prediction. For example, the *square footage* of a house times the *average housing price per square foot* in the neighborhood may be a good predictor of the *market value* of the house. Pretrained on large volumes of data, large language models (LLMs) (Radford et al., 2019; Brown et al., 2020; OpenAI, 2023; Touvron et al., 2023a;b) handle general language processing tasks and encapsulate extensive domain knowledge. Under proper instructions, an LLM can process dataset semantic information and utilize its knowledge to au-



**Figure 1:** Overview of PromptFE: (1) instructing the LLM to construct new features by providing dataset descriptions and example features; (2) evaluating the constructed features; (3) updating the prompt with top-performing features and scores; and (4) selecting a set of constructed features to add to the dataset.

054 tomatically construct features in a manner similar to domain experts. The work by Hollmann et al.  
055 (2023) demonstrates the potential of such research direction but is not sufficiently effective in fea-  
056 ture search. Similarly, the work by Nam et al. (2024) suffers from large search space. The works by  
057 Han et al. (2024) and Zhang et al. (2024) do not involve feature learning and improvement.

058 We present AutoFE by **Prompting** (PromptFE), a novel AutoFE framework that leverages LLMs for  
059 effective, efficient, and interpretable feature engineering, as illustrated in Figure 1. With dataset de-  
060 scriptions and example features in canonical Reverse Polish Notation (cRPN), we prompt the LLM  
061 to construct new features. After evaluation, we update the prompt with top-performing features with  
062 the evaluation scores and instruct the LLM to construct further features. Iteratively, the LLM ex-  
063 plores the feature space and improves solutions by learning good examples in context. The dataset  
064 semantic information not only guides feature search, but helps the LLM understand the patterns in  
065 example features. Applying domain knowledge, the LLM generates semantically meaningful fea-  
066 tures and explains their usefulness. Experiments on real-world datasets demonstrate that PromptFE  
067 yields over 5% mean performance gain for three downstream models and significantly outperforms  
068 state-of-the-art baselines. Furthermore, we show in ablation study the effects of dataset semantic  
069 context and proposed feature canonicalization scheme. We also comprehensively study the behavior  
070 of the LLM-based feature construction process.

071 Our main contributions are: (1) We introduce a novel LLM-based AutoFE framework utilizing  
072 dataset semantic information for automated feature construction, which is the first method capable  
073 of generating features in the RPN format while providing semantic explanations. (2) We benchmark  
074 the performance of our approach against state-of-the-art baselines using both GPT-3.5 and GPT-4.  
075 (3) We investigate the impact of semantic context and study the behavior of the LLM-based feature  
076 construction process, providing a comprehensive view of our approach.

## 077 2 RELATED WORK

080 **Large Language Models.** LLMs are large-scale general-purpose neural networks pretrained on vast  
081 corpora of text data, typically built with transformer-based architectures (Vaswani et al., 2017). Gen-  
082 erative LLMs, such as the GPT family (Radford et al., 2019; Brown et al., 2020; OpenAI, 2023) and  
083 the LLaMA family (Touvron et al., 2023a;b), are pretrained to successively generate the next token  
084 given the text input and can be finetuned using reinforcement learning from human feedback (Ziegler  
085 et al., 2019; Ouyang et al., 2022). By this means, they acquire the syntactic and semantic knowl-  
086 edge of natural languages and achieve state-of-the-art performance on various tasks including text  
087 generation, summarization, and question answering. Prompting techniques (Liu et al., 2023) have  
088 been developed to adapt LLMs to downstream tasks without modifying model weights. Few-shot  
089 learning (Brown et al., 2020) includes examples in the prompt for the language model to learn in  
090 context. Leveraging such capability, an LLM may function as a problem solver (Yang et al., 2024)  
091 that iteratively improves candidate solutions according to the task description and performance feed-  
092 back. Chain-of-thought (Wei et al., 2022; Kojima et al., 2022) strengthens reasoning performance  
093 of LLMs through the elicitation of intermediate reasoning steps.

094 **Automated Feature Engineering.** AutoFE complements the input dataset with engineered fea-  
095 tures to enhance the performance of downstream models. Traditional AutoFE approaches include  
096 expansion-reduction (Kanter & Veeramachaneni, 2015; Horn et al., 2020; Zhang et al., 2023), evo-  
097 lutionary algorithms (Smith & Bull, 2005; Zhu et al., 2022a), and reinforcement learning (Khu-  
098 rana et al., 2018; Li et al., 2023; Wang et al., 2023). DIFER (Zhu et al., 2022b) utilizes encoder-  
099 decoder neural networks to learn the utility of features and optimize features in the embedding  
100 space. OpenFE (Zhang et al., 2023) develops a feature boost algorithm to speedup feature eval-  
101 uation. Nonetheless, these traditional approaches do not incorporate the semantic information of  
102 datasets, which hampers the efficacy and interpretability of engineered features.

103 **AutoFE with Domain Knowledge.** The benefits of incorporating domain knowledge in AutoFE  
104 include: (1) improving the effectiveness; and (2) reducing the cost of feature search, especially  
105 the feature evaluation overhead. One direction in prior works is to learn transferrable knowledge.  
106 LFE (Nargesian et al., 2017) represents features with quantile sketches transferable across datasets  
107 and inputs them to a feature transformation recommendation model. FETCH (Li et al., 2023) is an  
RL-based AutoFE framework that takes tabular data as the state and is generalizable to new data. E-

AFE (Wang et al., 2023) pretrains a feature evaluator to help efficiently train the RL-based AutoFE model. The other direction is to leverage the semantic information of datasets. KAFE (Galhotra et al., 2019) employs knowledge graphs to identify semantically informative features relevant to the prediction task. CAAFE (Hollmann et al., 2023) manipulates datasets using the code generated from an LLM based on dataset descriptions. FeatLLM (Han et al., 2024) generates first-order rules for classification tasks. ELF-Gym (Zhang et al., 2024) generates first feature descriptions and then feature code. Neither approach involves feature learning and improvement. OCTree (Nam et al., 2024) relies on external decision tree algorithms to represent features and suffers from large search space. Differently, we adopt a compact form of feature representation in cRPN with pre-defined transformation operators. Our approach reduces the search space and helps the LLM learn the patterns of useful features, leading to stronger and more robust performance.

### 3 NOTATIONS

We denote a tabular dataset as  $D = \langle \mathbb{X}, \mathbf{y} \rangle$ , where  $\mathbb{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_d\}$  is the set of raw features with  $\mathbf{x}_i \in \mathbb{R}^n$  for  $i = 1, \dots, d$  and  $\mathbf{y} \in \mathbb{R}^n$  is the target. We construct a new feature  $\tilde{\mathbf{x}} = t(\mathbf{x}_{j_1}, \dots, \mathbf{x}_{j_o})$  by transforming existing features  $\mathbf{x}_{j_1}, \dots, \mathbf{x}_{j_o}$  via some operator  $t \in \mathbb{R}^n \times \dots \times \mathbb{R}^n \rightarrow \mathbb{R}^n$  of arity  $o$ . Given a set of transformation operators  $\mathbb{T}$ , we define the feature space  $\mathbb{X}_{\mathbb{T}}$  recursively as: for any  $\tilde{\mathbf{x}} \in \mathbb{X}_{\mathbb{T}}$ , either  $\tilde{\mathbf{x}} \in \mathbb{X}$ ; or  $\exists t \in \mathbb{T}$ , s.t.,  $\tilde{\mathbf{x}} = t(\tilde{\mathbf{x}}_{j_1}, \dots, \tilde{\mathbf{x}}_{j_o})$ , where  $\tilde{\mathbf{x}}_{j_1}, \dots, \tilde{\mathbf{x}}_{j_o} \in \mathbb{X}_{\mathbb{T}}$ . To measure feature complexity, we compute the order of a feature  $\tilde{\mathbf{x}} \in \mathbb{X}_{\mathbb{T}}$  as:

$$\alpha(\tilde{\mathbf{x}}) = \begin{cases} 0 & \text{if } \tilde{\mathbf{x}} \in \mathbb{X}, \\ 1 + \max_j \alpha(\tilde{\mathbf{x}}_j) & \text{if } \tilde{\mathbf{x}} = t(\tilde{\mathbf{x}}_{j_1}, \dots, \tilde{\mathbf{x}}_{j_o}) \exists t \in \mathbb{T}. \end{cases} \quad (1)$$

The constrained feature space with the order upper bounded by  $k$  is denoted as  $\mathbb{X}_{\mathbb{T}}^{(k)} = \{\tilde{\mathbf{x}} \in \mathbb{X}_{\mathbb{T}} \mid \alpha(\tilde{\mathbf{x}}) \leq k\}$ .

We denote the performance of a downstream machine learning model algorithm  $M$  on the dataset as  $\mathcal{E}_M(\mathbb{X}, \mathbf{y})$ . The objective of AutoFE is to augment the dataset with a set of constructed features  $\tilde{\mathbb{X}}^*$  to optimize the model performance, specifically:

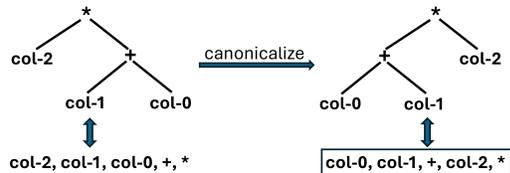
$$\tilde{\mathbb{X}}^* = \arg \max_{\tilde{\mathbb{X}} \subset \mathbb{X}_{\mathbb{T}}} \mathcal{E}_M(\mathbb{X} \cup \tilde{\mathbb{X}}, \mathbf{y}). \quad (2)$$

### 4 METHODOLOGY

In this section, we present PromptFE, a novel AutoFE framework leveraging the power of LLMs, particularly, the GPT models (Radford et al., 2019; Brown et al., 2020; OpenAI, 2023). The high-level idea is to provide the LLM with descriptive information of the dataset in the prompt and guide it to search for effective features using examples.

We represent features in a compact form in our prompt. A feature  $\tilde{\mathbf{x}} \in \mathbb{X}_{\mathbb{T}}$  is expressible as a tree, where the leaf nodes are raw features and the internal nodes are operators. However, the expression trees of features containing commutative operators (like addition and multiplication) are not unique since the child nodes of these operators are unordered. We introduce a canonicalization scheme: arranging operator nodes before feature nodes for left skewness and lexicographically sorting the nodes within each group. We then serialize the canonical expression tree into the postorder depth-first traversal string, i.e., canonical reverse Polish notation (cRPN), ensuring the one-to-one mapping between features and string representations. We denote the feature corresponding to an RPN string  $f$  as  $\tilde{\mathbf{x}}_f$  and the set of features corresponding to a set of RPN strings  $\mathbb{F}$  as  $\tilde{\mathbb{X}}_{\mathbb{F}}$ . We make further discussions in Appendix A.

Our prompt contains: (1) a meta description of the dataset; (2) an indexed list of the dataset attributes, with attribute types, value ranges, and descriptions; (3) lists of transformation operators with descriptions, grouped by the arity; (4) a ranked list of example features with performance evaluation scores; (5) feedback of previously constructed features; and (6) an output template of new



**Figure 2:** We obtain canonical RPN (cRPN) by re-ordering the nodes of a feature expression tree.

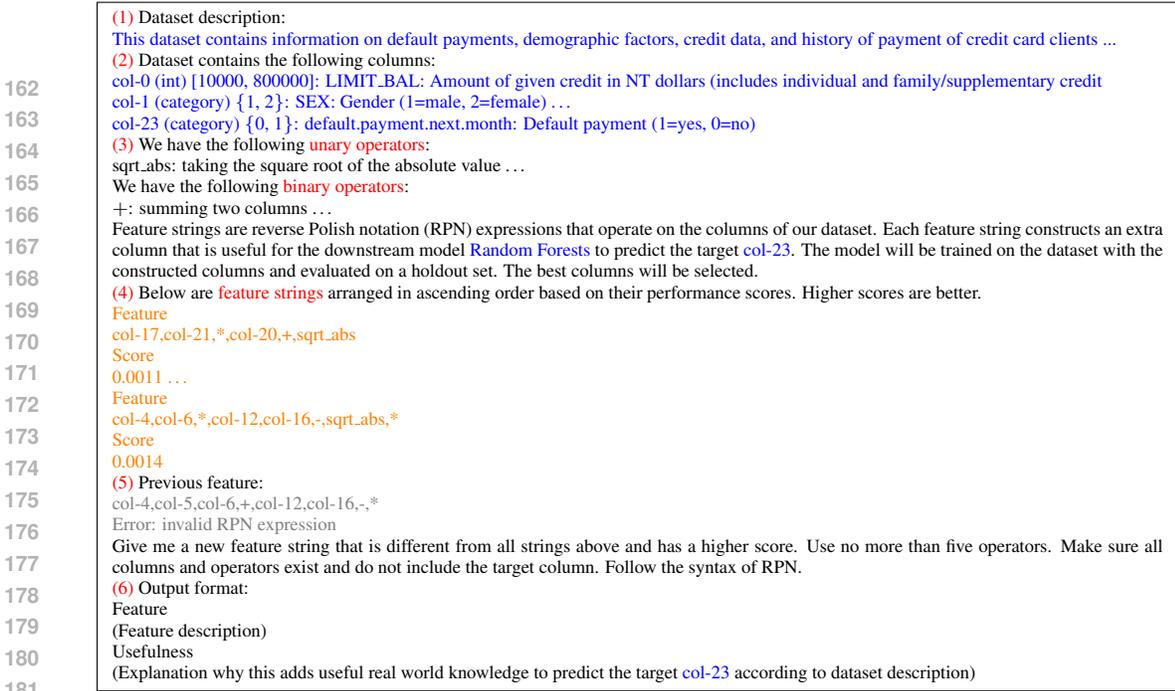


Figure 3: Prompt template. Sections containing dataset information are marked in blue. The ranked list of feature examples and scores is marked in orange. The feedback is marked in gray.

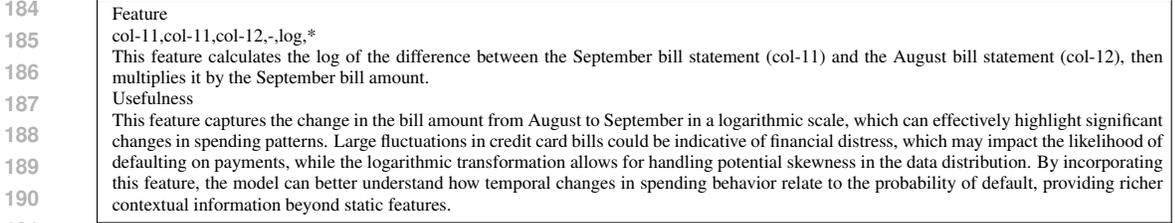


Figure 4: The LLM constructs a new feature in RPN and explains its usefulness from the semantic perspective.

features and explanations. Figure 3 outlines the structure of our prompt. The descriptions of the dataset and attributes provide contextual information for the LLM to understand the data and apply domain knowledge. The value ranges of attributes are useful for selecting appropriate feature transformations, e.g., min-max normalization when the scale is too large. We include the descriptions of transformation operators as they help the LLM parse example features in RPN syntax and construct syntactically valid feature strings. The output template not only structures the output but instructs the LLM to reason about the usefulness of the constructed features and offer semantic explanations, utilizing the chain-of-thought technique (Wei et al., 2022; Kojima et al., 2022). We additionally add a constraint instruction to use no more than a certain number of operators, which reduces the search space and regularizes the solutions. Figure 4 shows an example LLM output. The prompt may further include dataset statistics like mean, standard deviation, and skewness of the attributes.

We initialize the prompt with  $k$  random features from the constrained feature space  $\tilde{x}_1, \dots, \tilde{x}_k \in \mathbb{X}_{\mathbb{T}}^{(2)}$  represented in cRPN for demonstration, where the feature attributes are sampled per the softmax probabilities of feature importance by fitting the downstream model on the training data. This lets the LLM start search from a small feature space where it is easier to identify the basic patterns of promising features. Optionally, we can import external example features. We prompt the LLM to construct a fixed number of  $m$  new feature in an iteration. For each constructed feature string  $f$ , we first try to obtain the cRPN expression  $f^c$  to check whether  $f^c$  is syntactically valid and not a duplicate of candidate features. If both criteria are met, we evaluate the performance score of adding the single feature to the dataset  $s = \mathcal{E}_M(\mathbb{X} \cup \{\tilde{x}_{f^c}\}, \mathbf{y})$  through cross validation on the training data and add  $\langle f^c, s \rangle$  to the candidate set  $\mathbb{F}_{cand}$ . When  $f^c$  is among the top- $k$  candidate features in terms of the score  $s$ , we update prompt examples with the top- $k$  pairs  $\langle f', s' \rangle \in \mathbb{F}_{cand}$  ranked in the ascending order, taking score increment  $s' - \mathcal{E}_M(\mathbb{X}, \mathbf{y})$  from the baseline. We also provide the feedback of previously constructed features with scores or error messages for improvement. We then instruct the LLM to construct additional features using the updated prompt. To select candidate features,

**Algorithm 1: AutoFE by Prompting**


---

**Input :** Dataset  $D = (\mathbb{X}, \mathbf{y})$ , downstream model  $M$ , large language model  $LLM$ , and optionally an external set of features with evaluation scores  $\mathbb{F}_{ext}$

**Output:** A set of engineered features  $\mathbb{F}$

```

1 Initialize prompt  $P$  with dataset descriptions and example features;  $\mathbb{F}_{cand} \leftarrow \mathbb{F}_{ext}$  if  $\mathbb{F}_{ext}$  is available, otherwise  $\mathbb{F}_{cand} \leftarrow \emptyset$ ;  $\mathbb{F}_{set} \leftarrow \emptyset$ 
2 repeat
3    $\mathbb{F}_{LLM} = \{f_1, \dots, f_m\} \leftarrow LLM(P)$  ▷ Feature generation
4   for each  $f \in \mathbb{F}_{LLM}$  do
5      $f^c \leftarrow$  Canonicalize  $f$ 
6     if  $f^c$  is valid and  $f^c \notin \mathbb{F}_{cand}$  then ▷ Feature evaluation
7       Evaluate cross validation performance  $s \leftarrow \mathcal{E}_M(\mathbb{X} \cup \{x_{f^c}\}, \mathbf{y})$ 
8        $\mathbb{F}_{cand} \leftarrow \mathbb{F}_{cand} \cup \{f^c, s\}$ 
9     end
10  end
11  Update  $P$  such that  $P$  contains the top- $k$   $\langle f', s' \rangle \in \mathbb{F}_{cand}$  as ordered by  $s'$ 
12  if feature selection then
13    for  $n \leftarrow 1$  to  $|\mathbb{F}_{cand}|$  do ▷ Feature selection
14       $\mathbb{F}_n \leftarrow$  The top- $n$  features in  $\mathbb{F}_{cand}$  as ordered by  $s$ 
15      Evaluate validation performance  $s_n \leftarrow \mathcal{E}_M(\mathbb{X} \cup \mathbb{X}_{\mathbb{F}_n}, \mathbf{y})$ 
16    end
17     $\mathbb{F}_{set} \leftarrow \mathbb{F}_{set} \cup \{\langle \mathbb{F}_{n^*}, s_{n^*} \rangle\}$ , with  $n^* \leftarrow \operatorname{argmax}_n s_n$ 
18  end
19 until stopping criteria are met
20 return  $\mathbb{F}$  in  $\mathbb{F}_{set}$  with the maximum validation score

```

---

we successively add candidate features to the dataset from the best to the worst and determine the optimal number of features to add based on validation performance, which is evaluated over sets of candidate features and thus takes feature interactions into account.

Algorithm 1 summarizes our methodology. The size of the prompt scales linearly with the number of features in the dataset  $d$  and the number of example features  $k$  and stays roughly constant across feature construction iterations. Thus, the cost of an LLM generation step in line 3 is almost constant. The computation cost of feature evaluation in line 7 is also constant, preserving the efficiency and scalability of our algorithm. The evaluations in line 7 and at lines 13-16 are parallelizable.

In our algorithm, the LLM is instructed to perform as a problem solver (Yang et al., 2024). Analogous to evolutionary algorithms that generate new solutions through crossover and mutations on high-fitness candidates (Smith & Bull, 2005; Zhu et al., 2022a; Morris et al., 2024), we provide top-performing features in the prompt. By learning examples and scores in-context (Brown et al., 2020), the LLM recognizes the patterns of promising features and generates new features that are likely to be useful. It can make analogies to, modify, or combine example features in the prompt (Appendix F.3). Early in the search, we expect greater exploration due to the diversity of initial examples. As iterations progress, the LLM exploits promising feature spaces more, gradually refining the search until convergence. The dataset semantic information enhances the effectiveness of feature search through the guidance as a prior. The LLM’s temperature can be adjusted to balance exploration and exploitation, with higher temperatures encouraging more diverse solutions and lower temperatures favoring incremental changes to example features.

We adopt the same set of transformation operators  $\mathbb{T}$  as those in (Zhu et al., 2022b), including:

- Unary transformations: logarithm, reciprocal, square root, and min-max normalization;
- Binary transformations: addition, subtraction, multiplication, division, and modulo.

In min-max normalization, we take the statistics from the training data. Other transformations require only the information of a single instance. Hence, all transformations can be performed on an individual test instance without leaking other instances’ information. Data leakage (Overman et al., 2024) is an issue that has not been properly addressed in many existing AutoFE works.

## 5 EXPERIMENTS

### 5.1 EXPERIMENTAL SETUP

We benchmark performance on public real-world datasets from Kaggle and UCI repositories covering different domains (Appendix D.1). The descriptive information of datasets and attributes is retrieved from the sources without further processing. The downstream models we evaluate include linear models (LASSO for regression tasks and logistic regression for classification tasks), Random

**Table 1:** Comparison of overall performance. For each compared method, the left and right columns show the performance without and with post AutoFE parameter tuning of downstream model algorithms, respectively. The best results are boldfaced, and the second best results are underlined.

| Model          | Raw             | DIFER          |                       | OpenFE          |                | CAAFF           |                |                 |                 | OCTree<br>GPT-4 |                 | PromptFE (ours)       |                              |                              |                              |
|----------------|-----------------|----------------|-----------------------|-----------------|----------------|-----------------|----------------|-----------------|-----------------|-----------------|-----------------|-----------------------|------------------------------|------------------------------|------------------------------|
|                |                 |                |                       |                 |                | GPT-3.5         |                | GPT-4           |                 |                 |                 | GPT-3.5               |                              | GPT-4                        |                              |
| Linear Model   | 0.5636<br>14.00 | 0.6248<br>9.17 | 0.6369<br>5.83        | 0.5871<br>10.58 | 0.5866<br>9.92 | 0.5946<br>10.00 | 0.5941<br>9.50 | 0.5945<br>10.50 | 0.5946<br>9.83  | 0.6038<br>8.25  | 0.6044<br>7.33  | 0.6485<br>5.08        | 0.6487<br>3.50               | <b>0.6532</b><br><u>3.33</u> | <u>0.6526</u><br><b>3.17</b> |
| Random Forests | 0.7252<br>12.71 | 0.7400<br>8.29 | <u>0.7411</u><br>5.86 | 0.7380<br>8.07  | 0.7376<br>8.64 | 0.7387<br>6.29  | 0.7378<br>7.50 | 0.7357<br>9.00  | 0.7352<br>10.79 | 0.7348<br>9.14  | 0.7346<br>10.79 | 0.7408<br><b>4.43</b> | <b>0.7412</b><br><u>4.57</u> | 0.7392<br>6.29               | 0.7393<br>7.64               |
| LightGBM       | 0.7364<br>10.43 | 0.7504<br>8.86 | 0.7531<br>6.29        | 0.7454<br>9.50  | 0.7476<br>9.14 | 0.7457<br>9.07  | 0.7461<br>8.14 | 0.7405<br>9.00  | 0.7457<br>8.50  | 0.7409<br>11.21 | 0.7403<br>12.21 | 0.7522<br>5.71        | <b>0.7558</b><br><u>3.57</u> | <u>0.7542</u><br><b>3.57</b> | 0.7538<br>4.79               |
| Mean           | 0.6806          | 0.7091         | 0.7140                | 0.6953          | 0.6958         | 0.6979          | 0.6976         | 0.6950          | 0.6967          | 0.6976          | 0.6975          | 0.7171                | <u>0.7185</u>                | <b>0.7187</b>                | 0.7183                       |
| Mean Rank      | 12.30           | 8.75           | 6.00                  | 9.33            | 9.20           | 8.38            | 8.33           | 9.45            | 9.70            | 9.60            | 10.25           | 5.08                  | <b>3.90</b>                  | <u>4.45</u>                  | 5.30                         |

Forests (Breiman, 2001), and LightGBM (Ke et al., 2017). For linear models, we target-encode categorical features and min-max scale all features. We tune downstream model algorithm parameters by randomized search prior to and post AutoFE, because the model may need reconfiguration to accommodate the added features. Data are randomly split into training (64%), validation (16%), and test (20%) sets. We evaluate regression performance with  $1 - (\text{relative absolute error})^1$  and classification performance with accuracy. A higher evaluation score indicates better performance.

We compare PromptFE with the following state-of-the-art AutoFE methods: (1) DIFER (Zhu et al., 2022b): A neural network-based method that optimizes features in the embedding space using LSTMs to encode and decode features; (2) OpenFE (Zhang et al., 2023): An expansion-reduction method that evaluates features up to a certain order using a feature boost algorithm; (3) CAAFE (Hollmann et al., 2023): An LLM-based method that produces Python code to manipulate datasets stored in Pandas data frames; (4) OCTree (Nam et al., 2024): An LLM-based method that generates rules to manipulate datasets and encodes features using decision tree algorithms.

We employ `gpt-3.5-turbo-0125`<sup>2</sup> and `gpt-4-0613`<sup>2</sup> as the LLMs. For PromptFE, we include  $k = 10$  example features in the prompt and set the temperature of LLMs to 1 based on validation. We instruct the LLM to construct  $m = 1$  feature in each generation step for the best control of feature generation. We perform feature selection each time 10 new candidate features are constructed and terminate the algorithm once we have 200 candidate features. Parameters of the baseline methods are initialized per the corresponding papers. We make five repeated runs.

## 5.2 PERFORMANCE COMPARISON

Table 1 compares the overall performance between PromptFE and the baseline methods. Full results are presented in Appendix D.5<sup>3</sup>. PromptFE attains the best mean performance score and the lowest mean rank for all three downstream models, yielding over 5% mean performance gain and over 15% gain for linear models. We observe the greatest gain for linear models because unlike Random Forests and LightGBM, they cannot learn non-linear relationships themselves. The performance margin between PromptFE and baselines other than DIFER is statistically significant with  $p < 0.01$  by Friedman-Nemenyi test. PromptFE consistently outperforms CAAFE and OCTree, showing the robustness of PromptFE that reduces the search space with pre-defined operators and represents features in compact cRPN. Post-AutoFE parameter tuning brings the greatest performance improvement to DIFER, as it adds the most features to datasets (Appendix D.9). Compared with DIFER evaluating over 2000 candidate features during feature search, PromptFE evaluates only 200 candidate features (Appendix D.10). The higher efficiency of PromptFE is brought by the construction of semantically meaningful and effective features with the guidance of dataset semantic information.

We note that in PromptFE, using GPT-4 yields better performance for linear models but slightly worse performance for Random Forests than GPT-3.5. We speculate this is because the stronger in-context learning capability of GPT-4 increases the tendency of overfitting example features. One way to address this is to include more example features in the prompt to fully leverage GPT-4’s enhanced in-context learning capability (Appendix D.8).

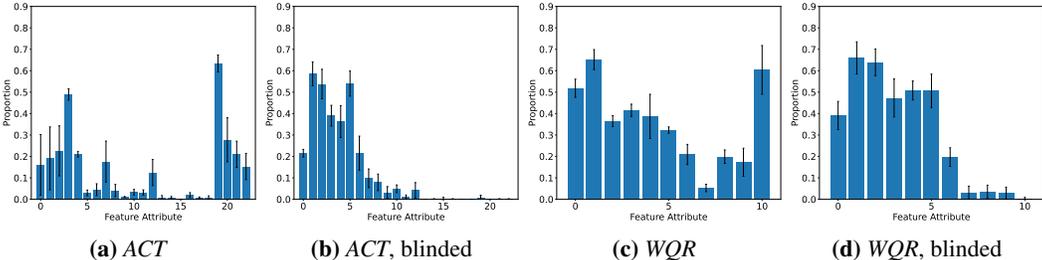
<sup>1</sup> $1 - \frac{\sum_i |y_i - \hat{y}_i|}{\sum_i |y_i - \bar{y}|}$ ,  $y$  is the target and  $\hat{y}$  is the prediction.

<sup>2</sup><https://platform.openai.com/docs/models>

<sup>3</sup>We were unable to complete the OCTree evaluations using GPT-3.5 as it easily got stuck with iteratively generating rules that triggered errors.

**Table 2:** Comparison of PromptFE with ablated versions. For each compared version, the left and middle columns show the performance without and with post AutoFE parameter tuning of downstream model algorithms, respectively, and the right column shows the number of LLM generations. Statistical significance of performance difference by Friedman-Nemenyi test is indicated with \* for  $p < 0.05$  and \*\* for  $p < 0.01$ .

|         | Model          | w/o Semantic Context |          |       | w/o Canonicalization |         |       | PromptFE |        |       |
|---------|----------------|----------------------|----------|-------|----------------------|---------|-------|----------|--------|-------|
| GPT-3.5 | Linear Model   | 0.6411               | 0.6433   | 443.4 | 0.6471               | 0.6486  | 349.1 | 0.6485   | 0.6487 | 356.7 |
|         | Random Forests | 0.7326**             | 0.7328** | 472.5 | 0.7372               | 0.7373  | 358.0 | 0.7408   | 0.7412 | 370.4 |
|         | LightGBM       | 0.7479*              | 0.7494   | 490.0 | 0.7485               | 0.7490  | 348.9 | 0.7522   | 0.7558 | 360.2 |
|         | Mean           | 0.7105**             | 0.7118** | 469.9 | 0.7141               | 0.7148  | 352.2 | 0.7171   | 0.7185 | 362.7 |
| GPT-4   | Linear Model   | 0.6437               | 0.6461   | 253.9 | 0.6462               | 0.6463  | 323.6 | 0.6532   | 0.6526 | 326.3 |
|         | Random Forests | 0.7285*              | 0.7288*  | 262.9 | 0.7366               | 0.7366  | 315.7 | 0.7392   | 0.7393 | 333.0 |
|         | LightGBM       | 0.7420**             | 0.7437   | 250.7 | 0.7461*              | 0.7480  | 328.5 | 0.7542   | 0.7538 | 335.7 |
|         | Mean           | 0.7078**             | 0.7092** | 255.9 | 0.7128**             | 0.7135* | 322.5 | 0.7187   | 0.7183 | 331.9 |



**Figure 6:** Distributions of feature attribute selection in the constructed features for linear models with GPT-4.

### 5.3 EFFECT OF SEMANTIC CONTEXT

We compare with the blinded version without dataset semantic information (Appendix C.2). From Table 2, PromptFE outperforms the blinded version for all downstream models with statistical significance. The performance difference is more pronounced for Random Forests and LightGBM, likely because the inclusion of non-semantically meaningful features by the blinded version consumes model capacity and causes greater overfitting to the training data. GPT-4 constructs features more efficiently than GPT-3.5 due to stronger capabilities. The incorporation of dataset semantic context improves the feature construction efficiency of GPT-3.5 but reduces that of GPT-4, as it guides to more focused feature spaces that increase the chances of duplication with candidate features.

### 5.4 EFFECT OF FEATURE EXPRESSION CANONICALIZATION

We compare with the ablated version without canonicalization of feature expressions. From Table 2, PromptFE outperforms the ablated version for all downstream models. Without canonicalization, we observe a slight decrease in the number of LLM generations. Since a feature can be represented in different expressions, the chances of duplication with the expressions of candidate features during feature search are reduced. However, the effectiveness of the features constructed by the LLM degrades in this setting due to increased difficulty in learning optimal feature patterns.

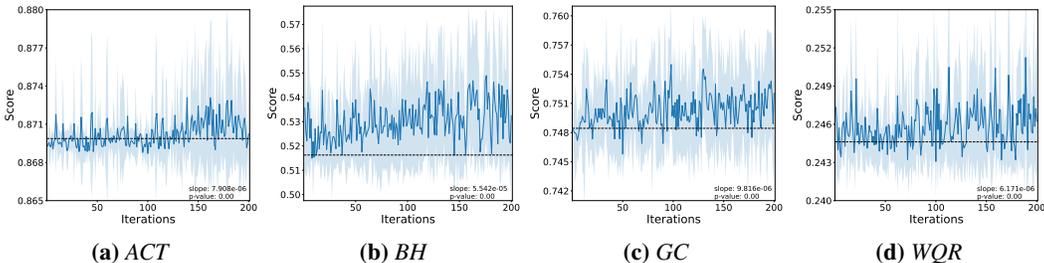
### 5.5 FEATURE ATTRIBUTE SELECTION

Figure 6 shows the distributions of feature attributes in the constructed features for linear models using GPT-4. Without semantic context, the LLM tends to prioritize earlier feature attributes in the dataset while paying less attention to later ones. In comparison, PromptFE is informed by the semantic context. Specifically, Attribute 19 *CD4 at baseline* in *ACT* and Attribute 10 *alcohol* in *WQR*, which contain crucial information for predicting the respective targets *censoring indicator* and *quality*, are consistently among the most frequent ones. This shows how the LLM leverages dataset semantic information to construct semantically meaningful and effective features in PromptFE.

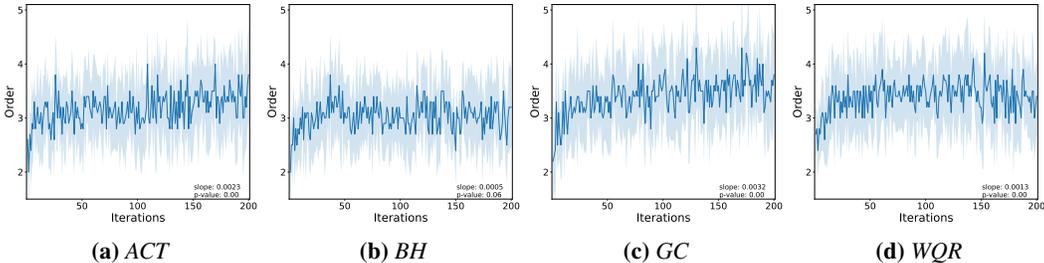
### 5.6 PERFORMANCE ANALYSIS

We study the performance for linear models with GPT-3.5 from ten repeated runs. Figures 7-10 display the slope and  $p$ -value from one-tailed t-tests in OLS regressions, with the shaded area showing one standard deviation above and below the mean curve.

378  
379  
380  
381  
382  
383  
384  
385  
386  
387  
388  
389  
390  
391  
392  
393  
394  
395  
396  
397  
398  
399  
400  
401  
402  
403  
404  
405  
406  
407  
408  
409  
410  
411  
412  
413  
414  
415  
416  
417  
418  
419  
420  
421  
422  
423  
424  
425  
426  
427  
428  
429  
430  
431



**Figure 7:** The cross validation score of candidate features on training data across iterations. The baseline score with raw features is indicated with the dash line.



**Figure 8:** The order of candidate features across iterations.

**Feature Learning.** We examine the cross validation score of candidate features across iterations. Figure 7 shows a significantly upward trend in the score, with most constructed features improving the performance. This demonstrates that PromptFE effectively improves the quality of constructed features through in-context learning of top-performing examples during feature search.

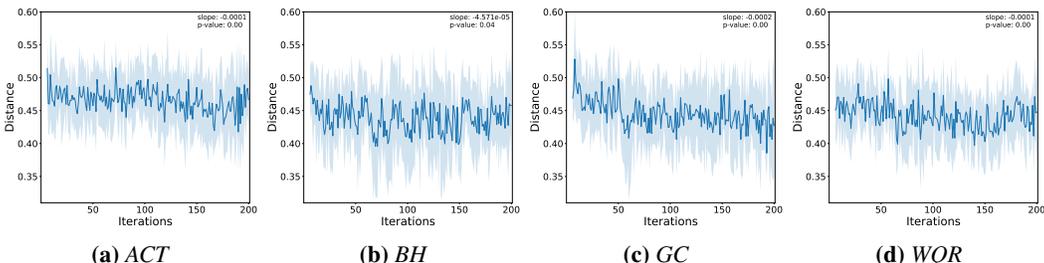
**Feature Complexity.** We examine the order of candidate features across iterations. Figure 8 shows that the feature order increases rapidly in early iterations and stabilizes over time. PromptFE effectively constructs complex features within promising feature spaces. Moreover, our constraint instruction offers regularization that prevents the construction of overly complex features.

**Feature Divergence.** We analyze the divergence of a new candidate feature from previous ones during feature search. We compute the edit distance between canonical feature expression trees using the algorithm by Zhang & Shasha (1989) and normalize the distance by the total number of nodes in both trees. Figure 9 shows the mean normalized tree edit distance between the current candidate feature and the previous five features across iterations. The observed downward trend indicates that feature search converges over iterations.

**Feature Construction Efficiency.** We examine the number of LLM generations needed to construct new candidate features across iterations. Figure 10 shows a slightly upward trend in the number of LLM generations, due to increasing difficulty of constructing non-duplicate features and higher likelihood of producing syntactical errors as features become more complex. Since the increase is non-significant, PromptFE remains scalable to a large number of iterations.

5.7 HYPERPARAMETER EFFECT

**Number of Examples in Prompt.** Table 3 reports the maximum validation score across iterations along with the number of LLM generations by varying the number of example features provided



**Figure 9:** The mean normalized edit distance between a candidate feature and previous five features.

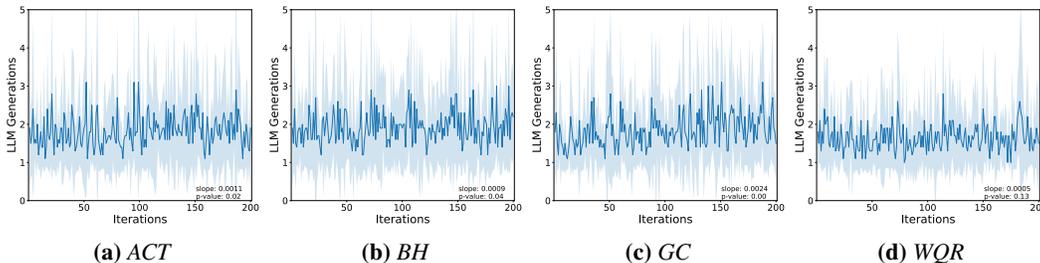


Figure 10: The number of LLM generations to construct a new candidate feature across iterations.

Table 3: Effect of the number of example features in the prompt with GPT-3.5. For each compared setting, the left column shows the validation score, and the right column shows the number of LLM generations.

| Model | Dataset | Number of Examples |       |               |       |               |              |               |              |
|-------|---------|--------------------|-------|---------------|-------|---------------|--------------|---------------|--------------|
|       |         | 1                  | 5     | 10            | 20    |               |              |               |              |
| RF    | AF      | 0.7895             | 507.2 | <b>0.7930</b> | 409.2 | 0.7914        | 393.2        | 0.7860        | <b>372.8</b> |
|       | WQR     | 0.3897             | 339.4 | 0.3937        | 329.8 | <b>0.3948</b> | 362.6        | 0.3940        | <b>330.0</b> |
|       | CD      | 0.8212             | 480.2 | 0.8213        | 371.2 | <b>0.8219</b> | 349.8        | 0.8218        | <b>343.2</b> |
| LGBM  | AF      | 0.8421             | 440.4 | <b>0.8433</b> | 404.6 | 0.8430        | <b>380.2</b> | 0.8420        | 384.2        |
|       | WQR     | 0.4248             | 346.0 | 0.4294        | 334.8 | 0.4301        | <b>322.8</b> | <b>0.4333</b> | 330.4        |
|       | CD      | <b>0.8228</b>      | 449.4 | 0.8224        | 361.2 | 0.8226        | 352.2        | 0.8228        | <b>321.2</b> |
| Mean  |         | 0.6817             | 427.1 | 0.6839        | 368.5 | <b>0.6840</b> | 360.1        | 0.6833        | <b>347.0</b> |

in the prompt. We observe that the best performance is achieved with 10 examples. Additionally, feature construction efficiency improves as the number of examples increases, as more examples can help the LLM reduce errors and generate more diverse features. Nonetheless, too many examples hinder the in-context learning of optimal feature patterns, as shown by the performance decline. The performance difference with 10 examples and with 1 example is statistically significant with  $p < 0.05$  by one-tailed paired t-tests.

**Temperature.** Table 4 reports the maximum validation score across iterations with the number of LLM generations under different LLM temperatures. We observe the best performance and efficiency when the temperature is around 1. Lower temperatures increase the likelihood of the LLM repeating previously constructed features, while higher temperatures made the LLM more prone to producing errors in the generations, both reducing feature construction efficiency. A temperature at 1 provides the best tradeoff between exploration and exploitation in feature search. The performance difference with the temperature at 1 and at 0.5 is statistically significant with  $p < 0.01$  by one-tailed paired t-tests.

Table 4: Effect of the LLM temperature with GPT-3.5.

| Model | Dataset | Temperature |        |               |              |               |       |
|-------|---------|-------------|--------|---------------|--------------|---------------|-------|
|       |         | 0.5         | 1      | 1.5           |              |               |       |
| RF    | AF      | 0.7875      | 794.4  | 0.7914        | <b>393.2</b> | <b>0.7916</b> | 609.2 |
|       | CD      | 0.8211      | 823.2  | <b>0.8219</b> | <b>349.8</b> | 0.8218        | 672.6 |
| LGBM  | AF      | 0.8365      | 1313.2 | <b>0.8430</b> | <b>380.2</b> | 0.8418        | 627.6 |
|       | CD      | 0.8225      | 519.8  | <b>0.8226</b> | <b>352.2</b> | 0.8223        | 662.6 |
| Mean  |         | 0.8169      | 862.7  | <b>0.8197</b> | <b>368.9</b> | 0.8194        | 643.0 |

## 6 CONCLUSION

In this paper, we present a novel LLM-based AutoFE framework for effective, efficient, and interpretable feature engineering that leverages the semantic information of datasets. It features an elegant approach to instructing the LLM to generate semantically meaningful features with explanations by providing dataset descriptions and example features in cRPN expressions. The LLM iteratively explores the feature space and improves feature construction by learning top-performing examples in context. We have demonstrated in extensive experiments that our approach significantly outperforms state-of-the-art AutoFE methods. The incorporation of semantic context from dataset descriptions and the proposed feature canonicalization scheme both contribute to performance improvement. We have also made comprehensive analysis on the LLM-based feature construction process. Our work opens up new possibilities for further LLM-driven applications on automated machine learning methodologies and underscores the potential of utilizing semantic information.

## 486 ETHICAL STATEMENT

487

488 All datasets used in this work are publicly available, free of personal information, and intended for  
489 research purposes only. Our use of GPT models complies with the terms and conditions of OpenAI.

490

## 491 REPRODUCIBILITY STATEMENT

492

493 The anonymized source code of this work can be accessed at [https://anonymous.4open.  
494 science/r/FEBP](https://anonymous.4open.science/r/FEBP).

495

## 496 REFERENCES

497

498 Leo Breiman. Random forests. *Machine learning*, 45:5–32, 2001.

499

500 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,  
501 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models  
502 are few-shot learners. In *Proceedings of Advances in neural information processing systems*,  
503 volume 33, pp. 1877–1901, 2020.

504 Milton Friedman. The use of ranks to avoid the assumption of normality implicit in the analysis of  
505 variance. *Journal of the american statistical association*, 32(200):675–701, 1937.

506

507 Sainyam Galhotra, Udayan Khurana, Oktie Hassanzadeh, Kavitha Srinivas, and Horst Samu-  
508 lowitz. Kafe: Automated feature enhancement for predictive modeling using external knowl-  
509 edge. In *Proceedings of NeurIPS 2019 Workshop: Knowledge Representation & Reasoning Meets  
Machine Learning*, 2019. URL [https://kr2ml.github.io/2019/papers/KR2ML\\_  
510 2019\\_paper\\_17.pdf](https://kr2ml.github.io/2019/papers/KR2ML_2019_paper_17.pdf).

511

512 Leo Grinsztajn, Edouard Oyallon, and Gael Varoquaux. Why do tree-based models still outperform  
513 deep learning on typical tabular data? In *Proceedings of the thirty-sixth Conference on Neural  
Information Processing Systems Datasets and Benchmarks Track*, 2022.

514

515 Sungwon Han, Jinsung Yoon, Sercan O Arik, and Tomas Pfister. Large language models can auto-  
516 matically engineer features for few-shot tabular learning. In *International Conference on Machine  
Learning*, pp. 17454–17479. PMLR, 2024.

517

518 Noah Hollmann, Samuel Müller, and Frank Hutter. Large language models for automated data  
519 science: Introducing CAAFE for context-aware automated feature engineering. In *Proceedings  
of the thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL [https:  
520 //openreview.net/forum?id=9WSxQZ9mG7](https://openreview.net/forum?id=9WSxQZ9mG7).

521

522 Franziska Horn, Robert Pack, and Michael Rieger. The autofeat python library for automated feature  
523 engineering and selection. In *Proceedings of Machine Learning and Knowledge Discovery in  
Databases*, pp. 111–120. Springer, 2020.

524

525 James Max Kanter and Kalyan Veeramachaneni. Deep feature synthesis: Towards automating data  
526 science endeavors. In *Proceedings of 2015 IEEE international conference on data science and  
advanced analytics (DSAA)*, pp. 1–10. IEEE, 2015.

527

528 Guolin Ke, Qi Meng, Thomas Finley, Taifeng Wang, Wei Chen, Weidong Ma, Qi-  
529 wei Ye, and Tie-Yan Liu. Lightgbm: A highly efficient gradient boosting decision  
530 tree. In *Proceedings of Advances in Neural Information Processing Systems*, volume 30,  
531 2017. URL [https://proceedings.neurips.cc/paper\\_files/paper/2017/  
532 file/6449f44a102fde848669bdd9eb6b76fa-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2017/file/6449f44a102fde848669bdd9eb6b76fa-Paper.pdf).

533

534 Udayan Khurana, Horst Samulowitz, and Deepak Turaga. Feature engineering for predictive mod-  
535 eling using reinforcement learning. In *Proceedings of the AAAI Conference on Artificial Intelli-  
536 gence*, volume 32, 2018.

537

538 Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large  
539 language models are zero-shot reasoners. In *Advances in Neural Information Processing Systems*,  
2022. URL <https://openreview.net/forum?id=e2TBb5y0yFf>.

- 540 Liyao Li, Haobo Wang, Liangyu Zha, Qingyi Huang, Sai Wu, Gang Chen, and Junbo Zhao. Learning  
541 a data-driven policy network for pre-training automated feature engineering. In *Proceedings of*  
542 *the eleventh International Conference on Learning Representations*, 2023.
- 543
- 544 Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. Pre-  
545 train, prompt, and predict: A systematic survey of prompting methods in natural language pro-  
546 cessing. *ACM Comput. Surv.*, 55(9), 2023. URL <https://doi.org/10.1145/3560815>.
- 547
- 548 Clint Morris, Michael Jurado, and Jason Zutty. Llm guided evolution - the automation of models  
549 advancing models. In *Proceedings of the Genetic and Evolutionary Computation Conference*,  
550 GECCO '24, pp. 377–384, New York, NY, USA, 2024. Association for Computing Machinery.  
551 URL <https://doi.org/10.1145/3638529.3654178>.
- 552 Jaehyun Nam, Kyuyoung Kim, Seunghyuk Oh, Jihoon Tack, Jaehyung Kim, and Jinwoo Shin. Op-  
553 timized feature generation for tabular data via llms with decision tree reasoning. *Advances in*  
554 *Neural Information Processing Systems*, 37:92352–92380, 2024.
- 555
- 556 Fatemeh Nargesian, Horst Samulowitz, Udayan Khurana, Elias B Khalil, and Deepak S Turaga.  
557 Learning feature engineering for classification. In *Proceedings of International Joint Conference*  
558 *on Artificial Intelligence*, volume 17, pp. 2529–2535, 2017.
- 559
- 560 Peter Bjorn Nemenyi. *Distribution-free multiple comparisons*. Princeton University, 1963.
- 561
- 562 OpenAI. Gpt-4 technical report, 2023. URL [https://cdn.openai.com/papers/gpt-4.](https://cdn.openai.com/papers/gpt-4.pdf)  
563 pdf.
- 564
- 565 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin,  
566 Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language  
567 models to follow instructions with human feedback. In *Proceedings of Advances in Neu-*  
568 *ral Information Processing Systems*, volume 35, pp. 27730–27744. Curran Associates, Inc.,  
569 2022. URL [https://proceedings.neurips.cc/paper\\_files/paper/2022/](https://proceedings.neurips.cc/paper_files/paper/2022/file/bl1efde53be364a73914f58805a001731-Paper-Conference.pdf)  
570 [file/bl1efde53be364a73914f58805a001731-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2022/file/bl1efde53be364a73914f58805a001731-Paper-Conference.pdf).
- 571
- 572 Tom Overman, Diego Klabjan, and Jean Utke. Iife: Interaction information based automated feature  
573 engineering, 2024. URL <https://arxiv.org/abs/2409.04665>.
- 574
- 575 Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Lan-  
576 guage models are unsupervised multitask learners, 2019. URL [https://d4mucfpxsywv.](https://d4mucfpxsywv.cloudfront.net/better-language-models/language-models.pdf)  
577 [cloudfront.net/better-language-models/language-models.pdf](https://d4mucfpxsywv.cloudfront.net/better-language-models/language-models.pdf).
- 578
- 579 Matthew G Smith and Larry Bull. Genetic programming with a genetic algorithm for feature con-  
580 struction and selection. *Genetic Programming and Evolvable Machines*, 6:265–281, 2005.
- 581
- 582 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée  
583 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Ar-  
584 mand Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation  
585 language models, 2023a. URL <https://arxiv.org/abs/2302.13971>.
- 586
- 587 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Niko-  
588 lay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruiti Bhosale, et al. Llama 2: Open founda-  
589 tion and fine-tuned chat models, 2023b. URL <https://arxiv.org/abs/2307.09288>.
- 590
- 591 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N  
592 Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Proceed-*  
593 *ings of Advances in Neural Information Processing Systems*, volume 30, 2017. URL  
[https://proceedings.neurips.cc/paper\\_files/paper/2017/file/](https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf)  
[3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf).
- 594
- 595 Kafeng Wang, Pengyang Wang, and Chengzhong Xu. Toward efficient automated feature engineer-  
596 ing. In *Proceedings of 2023 IEEE 39th International Conference on Data Engineering (ICDE)*,  
597 pp. 1625–1637, 2023. doi: 10.1109/ICDE55515.2023.00128.

- 594 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed H. Chi,  
595 Quoc V Le, and Denny Zhou. Chain of thought prompting elicits reasoning in large lan-  
596 guage models. In *Advances in Neural Information Processing Systems*, 2022. URL [https://openreview.net/forum?id=\\_VjQlMeSB\\_J](https://openreview.net/forum?id=_VjQlMeSB_J).  
597
- 598 Chengrun Yang, Xuezhi Wang, Yifeng Lu, Hanxiao Liu, Quoc V Le, Denny Zhou, and Xinyun  
599 Chen. Large language models as optimizers. In *Proceedings of the twelfth International Confer-  
600 ence on Learning Representations*, 2024. URL <https://openreview.net/forum?id=Bb4VGOWELI>.  
601
- 602 Kaizhong Zhang and Dennis Shasha. Simple fast algorithms for the editing distance between trees  
603 and related problems. *SIAM J. Comput.*, 18:1245–1262, 12 1989. doi: 10.1137/0218082.  
604
- 605 Tianping Zhang, Zheyu Zhang, Zhiyuan Fan, Haoyan Luo, Fengyuan Liu, Qian Liu, Wei Cao, and  
606 Jian Li. Openfe: automated feature generation with expert-level performance. In *Proceedings of  
607 the 40th International Conference on Machine Learning, ICML’23*, 2023.  
608
- 609 Yanlin Zhang, Ning Li, Quan Gan, Weinan Zhang, David Wipf, and Minjie Wang. Elf-gym: Evalu-  
610 ating large language models generated features for tabular prediction. In *Proceedings of the 33rd  
611 ACM International Conference on Information and Knowledge Management*, pp. 5420–5424,  
612 2024.
- 613 Guanghui Zhu, Shen Jiang, Xu Guo, Chunfeng Yuan, and Yihua Huang. Evolutionary automated  
614 feature engineering. In *Proceedings of Pacific Rim International Conference on Artificial Intelli-  
615 gence*, pp. 574–586. Springer, 2022a.  
616
- 617 Guanghui Zhu, Zhuoer Xu, Chunfeng Yuan, and Yihua Huang. Difer: Differentiable automated  
618 feature engineering. In *Proceedings of the first International Conference on Automated Machine  
619 Learning*, volume 188 of *Proceedings of Machine Learning Research*, pp. 17/1–17. PMLR, 25–27  
620 Jul 2022b. URL <https://proceedings.mlr.press/v188/zhu22a.html>.
- 621 Daniel M Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B Brown, Alec Radford, Dario Amodei, Paul  
622 Christiano, and Geoffrey Irving. Fine-tuning language models from human preferences, 2019.  
623 URL <https://arxiv.org/abs/1909.08593>.  
624

## 625 A DISCUSSION ON CANONICAL RPN FEATURE REPRESENTATION

### 626 A.1 WHY RPN

627 RPN provides a compact and unambiguous form of feature representation. In contrast, infix expres-  
628 sion requires extra information such as brackets to determine operator precedence. Without brackets,  
629 the feature in infix expression  $col-0 - (col-1 + col-2)$  would be indistinguishable from the feature  
630  $(col-0 - col-1) + col-2$ , while both features are distinctively encoded in RPN. Such compactness  
631 and unambiguity of RPN facilitate sequential modeling since there is no need to model the extra  
632 information, e.g., the positions of brackets.  
633

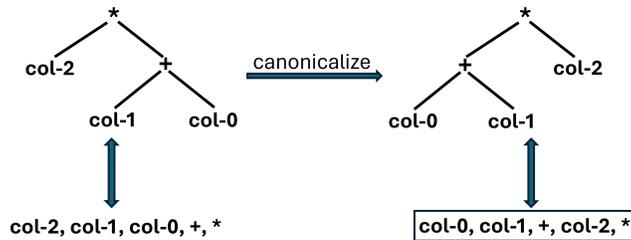
634 Compared with other forms of feature representation such as prefix expression of depth-first traversal  
635 or breadth-first traversal, RPN better encodes the recursive structure of the expression tree. The  
636 bottom-up enumeration of tree nodes makes it easy for the LLM to evaluate the feature expression  
637 by scanning the sequence from left to right, for instance,  $((col-0 col-1 -) col-2 +)$  (parentheses  
638 denote recursion). Using the prefix expression  $(+ (- col-0 col-1) col-2)$  or breadth-first expression  
639  $(+ (- [col-2] col-0 col-1))$ , however, the LLM always needs to look back to find the operator, which  
640 undermines sequential modeling. We find in our experiments that when using prefix expression, the  
641 LLM encounters difficulty in generating syntactically valid feature expressions.  
642

### 643 A.2 WHY CANONICALIZATION

644 Although there is one-to-one mapping between feature expression trees and RPN expressions, a fea-  
645 ture that contains commutative operators (like addition and multiplication) can be represented by  
646 different RPN expressions, since the child nodes of these operators are unordered. We introduce a  
647

648 canonicalization scheme: arranging operator nodes before feature nodes and lexicographically sort-  
 649 ing the nodes within each group. Through canonicalization, we create one-to-one mapping between  
 650 features and cRPN expressions. This ensures the consistency of our feature representations and  
 651 facilitates the in-context learning of feature patterns.

652 By arranging operator nodes before feature nodes, we also introduce left skewness to the expression  
 653 tree that enhances the clarity of the recursive structure in cRPN. As illustrated in Figure 11, the  
 654 original feature expression (col-2 (col-1 col-0 +) \*) becomes ((col-0 col-1 +) col-2 \*) after canoni-  
 655 calization, so that the LLM does not need to look back for col-2 when evaluating the expression.  
 656



665  
 666 **Figure 11:** Our canonicalization scheme introduces left skewness to the expression tree.

## 667 B CONVERSION BETWEEN FEATURE EXPRESSION TREE AND RPN

671 Algorithms 2 and 3 detail the process of conversion between a feature expression tree and an RPN  
 672 feature string. We check the RPN syntactical validity of a feature string in Algorithm 3 by checking  
 673 whether there is enough child node in the stack in line 6 and the size of the stack is exactly one (the  
 674 root) in line 13 returning the output.

---

### 676 **Algorithm 2:** Feature Expression Tree to RPN

---

677 **Input** : A feature expression tree  $T$

678 **Output:** An RPN feature string  $f$

```

679 1  $r \leftarrow$  the root of  $T$ 
680 2 Initialize string  $f \leftarrow \epsilon$ , stack  $S \leftarrow [r]$ , and  $visited \leftarrow \emptyset$ 
681 3 repeat
682 4    $u \leftarrow S.peek()$ 
683 5   if  $u \in visited$  then
684 6      $f.append(u)$ 
685 7      $S.pop()$ 
686 8   end
687 9   else
688 10    for each child  $v$  of  $u$  in the reverse order do
689 11       $S.push(v)$ 
690 12    end
691 13     $visited \leftarrow visited \cup \{u\}$ 
692 14  end
693 15 until  $S$  is empty
694 16 return  $f$ 

```

---

## 695 C EXAMPLE PROMPT

### 696 C.1 FULL PROMPT

697 Figure 12 shows an example of full prompts used in our main experiments.  
 698  
 699  
 700  
 701

**Algorithm 3:** RPN to Feature Expression Tree

---

```

702 Input : An RPN feature string  $f$ 
703 Output: The root of a feature expression tree  $T$ 
704
705 1 Initialize stack  $S \leftarrow []$ 
706 2 for  $i \leftarrow 1$  to  $|f|$  do
707 3    $u \leftarrow$  the  $i$ -th element of  $f$ 
708 4   if  $u$  is an operator then
709 5      $o \leftarrow$  the arity of  $u$ 
710 6     for  $j \leftarrow 1$  to  $o$  do
711 7        $v \leftarrow S.pop()$ 
712 8       Prepend  $v$  to the list of children of  $u$ 
713 9     end
714 10   end
715 11    $S.push(u)$ 
716 12 end
717 13 return  $S.pop()$ 

```

---

Figure 12: Example full prompt on the Credit Default dataset.

**Dataset description:**

This dataset contains information on default payments, demographic factors, credit data, history of payment, and bill statements of credit card clients in Taiwan from April 2005 to September 2005.

Dataset contains the following **columns**:

col-0 (int) [10000, 800000]: LIMIT\_BAL: Amount of given credit in NT dollars (includes individual and family/supplementary credit)

col-1 (category) {1, 2}: SEX: Gender (1=male, 2=female)

col-2 (category) {0, 1, 2, 3, 4, 5, 6}: EDUCATION: (1=graduate school, 2=university, 3=high school, 4=others, 5=unknown, 6=unknown)

col-3 (category) {0, 1, 2, 3}: MARRIAGE: Marital status (1=married, 2=single, 3=others)

col-4 (int) [21, 79]: AGE: Age in years

col-5 (category) {-2, -1, 0, 1, 2, 3, 4, 5, 6, 7, 8}: PAY\_0: Repayment status in September, 2005 (-1=pay dully, 1=payment delay for one month, 2=payment delay for two months, ... 8=payment delay for eight months, 9=payment delay for nine months and above)

...

col-23 (category) {0, 1}: default.payment.next.month: Default payment (1=yes, 0=no)

We have the following **unary operators**:

log: taking the log of the absolute value

sqrt\_abs: taking the square root of the absolute value

min\_max: min-max normalization

reciprocal: taking the reciprocal

We have the following **binary operators**:

+: summing two columns

-: subtracting two columns

\*: multiplying two columns

/: taking the division of two columns

mod\_column: taking the modulo of two columns

Feature strings are reverse Polish notation (RPN) expressions that operate on the columns of our dataset.

Each feature string constructs an extra column that is useful for the downstream model Random Forests to predict the target col-23. The model will be trained on the dataset with the constructed columns and evaluated on a holdout set. The best columns will be selected.

Below are **feature strings** arranged in ascending order based on their performance scores. Higher scores are better.

Feature

col-17,col-21,\*,col-20,+,sqrt\_abs

Score

0.0011

...

Feature

col-4,col-6,\*,col-12,col-16,-,sqrt\_abs,\*

756 Score  
 757 0.0014  
 758  
 759 Previous feature:  
 760 col-4,col-5,col-6+,col-12,col-16,-,\*  
 761 Error: invalid RPN expression  
 762  
 763 Give me a new feature string that is different from all strings above and has a higher score. Use  
 764 no more than five operators. Make sure all columns and operators exist and do not include the target  
 765 column. Follow the syntax of RPN.  
 766 **Output format:**  
 767 Feature  
 768 (Feature name and description)  
 769  
 770 Usefulness  
 771 (Explanation why this adds useful real world knowledge to predict the target col-23 according to dataset  
 772 description)  
 773

## 774 C.2 SEMANTICALLY BLINDED PROMPT

775  
 776 Figure 13 shows an example of semantically blinded prompts used in our experiments in Section 5.3.

777 Figure 13: Example semantically blinded prompt on the Credit Default dataset.

779 Dataset contains the following **columns**:  
 780 col-0  
 781 col-1  
 782 col-2  
 783 col-3  
 784 col-4  
 785 col-5  
 786 ...  
 787 col-23  
 788 We have the following **unary operators**:  
 789 log: taking the log of the absolute value  
 790 sqrt\_abs: taking the square root of the absolute value  
 791 min\_max: min-max normalization  
 792 reciprocal: taking the reciprocal  
 793 We have the following **binary operators**:  
 794 +: summing two columns  
 795 -: subtracting two columns  
 796 \*: multiplying two columns  
 797 /: taking the division of two columns  
 798 mod\_column: taking the modulo of two columns  
 799 Feature strings are reverse Polish notation (RPN) expressions that operate on the columns of our dataset.  
 800 Each feature string constructs an extra column that is useful for the downstream model Random Forests  
 801 to predict the target col-23. The model will be trained on the dataset with the constructed columns and  
 802 evaluated on a holdout set. The best columns will be selected.  
 803 Below are **feature strings** arranged in ascending order based on their performance scores. Higher scores  
 804 are better.  
 805  
 806 Feature  
 807 col-17,col-21,\*,col-20+,sqrt\_abs  
 808 Score  
 809 0.0011  
 810 ...  
 811 Feature  
 812 col-4,col-6,\*,col-12,col-16,-,sqrt\_abs,\*  
 813 Score  
 814 0.0014

810  
811  
812  
813  
814  
815  
816  
817  
818  
819  
820  
821  
822  
823  
824  
825  
826  
827  
828  
829  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
860  
861  
862  
863

Previous feature:  
col-4,col-5,col-6+,col-12,col-16,-,\*  
Error: invalid RPN expression

Give me a new feature string that is different from all strings above and has a higher score. Use no more than five operators. Make sure all columns and operators exist and do not include the target column. Follow the syntax of RPN.

**Output format:**  
Feature  
(Feature name and description)  
Usefulness  
(Explanation why this adds useful real world knowledge to predict the target col-23 according to dataset description)

**Table 5:** Statistics of datasets. The datasets cover different domains and tasks and vary in sizes.

| Name                       | Task           | # Samples | # Features | # Numerical | # Categorical |
|----------------------------|----------------|-----------|------------|-------------|---------------|
| Airfoil (AF)               | Regression     | 1,503     | 5          | 5           | 0             |
| Boston Housing (BH)        | Regression     | 506       | 13         | 12          | 1             |
| Bikeshare (BS)             | Regression     | 731       | 10         | 6           | 4             |
| Wine Quality Red (WQR)     | Regression     | 1,599     | 11         | 11          | 0             |
| AIDS Clinical Trials (ACT) | Classification | 2,139     | 23         | 9           | 14            |
| Credit Default (CD)        | Classification | 30,000    | 23         | 14          | 9             |
| German Credit (GC)         | Classification | 1,000     | 20         | 10          | 10            |

## D EXPERIMENTAL DETAIL

### D.1 DATASETS

Table 6 summarizes the sources of datasets used in our experiments. Datasets are selected such that they cover different domains and both regression and classification tasks. Most of them have been used in previous works (Zhu et al., 2022a;b; Zhang et al., 2023; Hollmann et al., 2023).

**Table 6:** Sources of datasets.

| Name                       | Source  |
|----------------------------|---|
| Airfoil (AF)               | <a href="https://archive.ics.uci.edu/dataset/291/airfoil+self+noise">https://archive.ics.uci.edu/dataset/291/airfoil+self+noise</a>                                     |
| Boston Housing (BH)        | <a href="https://www.kaggle.com/datasets/arunjangir245/boston-housing-dataset">https://www.kaggle.com/datasets/arunjangir245/boston-housing-dataset</a>                 |
| Bikeshare (BS)             | <a href="https://www.kaggle.com/datasets/marklvl/bike-sharing-dataset">https://www.kaggle.com/datasets/marklvl/bike-sharing-dataset</a>                                 |
| Wine Quality Red (WQR)     | <a href="https://archive.ics.uci.edu/dataset/186/wine+quality">https://archive.ics.uci.edu/dataset/186/wine+quality</a>   |
| AIDS Clinical Trials (ACT) | <a href="https://archive.ics.uci.edu/dataset/890/aids+clinical+trials+group+study+175">https://archive.ics.uci.edu/dataset/890/aids+clinical+trials+group+study+175</a> |
| Credit Default (CD)        | <a href="https://www.kaggle.com/datasets/uciml/default-of-credit-card-clients-dataset">https://www.kaggle.com/datasets/uciml/default-of-credit-card-clients-dataset</a> |
| German Credit (GC)         | <a href="https://archive.ics.uci.edu/dataset/573/south+german+credit+update">https://archive.ics.uci.edu/dataset/573/south+german+credit+update</a>                     |

### D.2 EXPERIMENTAL PLATFORM

All experiments are conducted on the Ubuntu 22.04.4 LTS operating system, 16 Intel(R) Core(TM) i7-7820X CPUs, and 4 NVIDIA GeForce RTX 2080 Ti GPUs, with the framework of Python 3.11.9 and PyTorch 1.12.1.

### D.3 FEATURE TRANSFORMATION OPERATORS

We list the details of all feature transformation operators below.

Unary transformations:

- Logarithm: Element-wise logarithm of the absolute value;
- Reciprocal: Element-wise reciprocal;
- Square root: Element-wise square root of the absolute value;
- Min-max normalization: Element-wise min-max normalization using the min and max values from the training data.

Binary transformations:

- Addition: Element-wise addition;
- Subtraction: Element-wise subtraction;
- Multiplication: Element-wise multiplication;
- Division: Element-wise division;
- Modulo: Element-wise modulo.

#### D.4 PARAMETER TUNING OF DOWNSTREAM MODELS

We tune the parameters of downstream models prior to and post AutoFE using randomized search implemented in an `Sklearn` package<sup>4</sup>. Table 7 lists the configurations of parameter tuning for each downstream model. We set the number of randomized search iterations to 100.

**Table 7:** Hyperparameter search space for downstream models.

| Model          | Parameter        | Search Space*            |
|----------------|------------------|--------------------------|
| Linear Model   | regularization   | loguniform(0.00001, 100) |
| Random Forests | num estimators   | randint(5, 250)          |
|                | max depth        | randint(1, 250)          |
|                | max features     | uniform(0.01, 0.99)      |
|                | max samples      | uniform(0.1, 0.9)        |
| LightGBM       | num estimators   | randint(10, 1000)        |
|                | num leaves       | randint(8, 64)           |
|                | learning rate    | loguniform(0.001, 1)     |
|                | bagging fraction | uniform(0.1, 0.9)        |
|                | feature fraction | uniform(0.1, 0.9)        |
|                | reg lambda       | loguniform(0.001, 100)   |

\* As defined in the `scipy.stats` package <https://docs.scipy.org/doc/scipy/reference/stats.html>.

#### D.5 FULL RESULTS

Tables 8-10 provide the full experimental results corresponding to the results in Tables 1 and 2. Tables 11-14 report the sample standard deviations corresponding to the experimental results in Tables 1-3, respectively.

#### D.6 RELATIVE PERFORMANCE IMPROVEMENT

Tables 15 and 16 report the percentage performance improvement of FEBP over the baseline methods with GPT-3.5 and GPT-4, respectively, corresponding to the experimental results in Table 1.

#### D.7 STATISTICAL TESTS

We perform the Friedman test (Friedman, 1937) to determine whether there is statistically significant difference among the compared AutoFE methods. The Friedman test  $p$ -values for the results in Tables 1 and 2 are  $1.16 \times 10^{-47}$  and  $3.95 \times 10^{-34}$ , respectively. Hence, we can reject the null hypothesis that the performance is the same for all methods. We perform the Nemenyi post-hoc test (Nemenyi, 1963) to further determine which AutoFE methods have different performance. Tables 17 and 18 summarize the  $p$ -values for the pairwise comparisons in Tables 1 and 2, respectively. From Table 17, the performance difference between our method FEBP and baseline methods other than DIFER (Zhu et al., 2022b) is statistically significant at the  $p = 0.01$  level. From Table 18, the performance difference between the full version of FEBP and the semantically blinded version is statistically significant at the  $p = 0.01$  level.

To examine the performance difference when using Random Forests and LightGBM, we perform additional statistical tests for the results in Table 1 excluding the linear model results. The Friedman test  $p$ -value is  $6.14 \times 10^{-23}$ . Table 19 summarizes the  $p$ -values from the Nemenyi post-hoc test for pairwise comparison. We observe that FEBP with GPT-3.5 and post-AutoFE parameter tuning significantly outperforms all baselines other than DIFER at the  $p = 0.05$  level. With GPT-4, the performance difference between FEBP and CAAFE (Hollmann et al., 2023) is statistically significant at the  $p = 0.05$  level.

<sup>4</sup>[https://scikit-learn.org/1.5/modules/generated/sklearn.model\\_selection.RandomizedSearchCV.html](https://scikit-learn.org/1.5/modules/generated/sklearn.model_selection.RandomizedSearchCV.html)

**Table 8:** Summary of experimental results. For each compared method, the left and right columns show the results without and with parameter tuning of the downstream model algorithm post AutoFE, respectively. The best results are highlighted in boldface, and the second best results are underlined.

| Model          | Dataset | Raw    | DIFER         |               | OpenFE        |               | CAAFE         |               |               |               | PromptFE (ours) |               |               |               |
|----------------|---------|--------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|-----------------|---------------|---------------|---------------|
|                |         |        |               |               |               |               | GPT-3.5       |               | GPT-4         |               | GPT-3.5         |               | GPT-4         |               |
| Linear Model   | AF      | 0.3474 | 0.5870        | 0.6090        | 0.4300        | 0.4303        | 0.4011        | 0.4016        | 0.4376        | 0.4378        | 0.6612          | 0.6616        | <b>0.6649</b> | <u>0.6647</u> |
|                | BH      | 0.3776 | 0.5013        | 0.4994        | 0.3900        | 0.3880        | 0.4788        | 0.4765        | 0.4503        | 0.4506        | 0.4995          | 0.5025        | <u>0.5184</u> | <b>0.5289</b> |
|                | WQR     | 0.2696 | 0.2475        | 0.2630        | 0.2713        | 0.2736        | 0.2742        | 0.2757        | <b>0.2776</b> | <u>0.2776</u> | 0.2722          | 0.2745        | 0.2713        | 0.2748        |
|                | ACT     | 0.8505 | 0.8715        | <b>0.8799</b> | 0.8729        | 0.8729        | 0.8519        | 0.8514        | 0.8565        | 0.8570        | 0.8729          | <u>0.8794</u> | 0.8766        | 0.8762        |
|                | CD      | 0.8267 | 0.8273        | 0.8280        | 0.8265        | 0.8268        | 0.8265        | 0.8267        | 0.8238        | 0.8238        | 0.8282          | 0.8282        | 0.8288        | <b>0.8288</b> |
|                | GC      | 0.7100 | 0.7140        | 0.7420        | 0.7320        | 0.7280        | 0.7350        | 0.7330        | 0.7210        | 0.7210        | <u>0.7570</u>   | 0.7460        | <b>0.7590</b> | 0.7420        |
| Mean           |         | 0.5636 | 0.6248        | 0.6369        | 0.5871        | 0.5866        | 0.5946        | 0.5941        | 0.5945        | 0.5946        | 0.6485          | 0.6487        | <b>0.6532</b> | <u>0.6526</u> |
| Mean Rank      |         | 12.00  | 8.17          | 5.50          | 9.25          | 8.50          | 8.67          | 8.17          | 8.83          | 8.17          | 4.75            | 3.17          | <u>3.00</u>   | <b>2.83</b>   |
| Random Forests | AF      | 0.7677 | 0.7650        | <u>0.7786</u> | 0.7579        | 0.7682        | 0.7711        | 0.7693        | 0.7696        | 0.7720        | 0.7709          | <b>0.7787</b> | 0.7681        | 0.7749        |
|                | BH      | 0.5426 | <b>0.5718</b> | <u>0.5701</u> | 0.5658        | 0.5620        | 0.5556        | 0.5556        | 0.5512        | 0.5492        | 0.5549          | 0.5533        | 0.5543        | 0.5522        |
|                | BS      | 0.9446 | 0.9865        | 0.9871        | 0.9901        | <u>0.9901</u> | <b>0.9916</b> | <b>0.9916</b> | 0.9818        | 0.9816        | 0.9873          | 0.9881        | 0.9845        | 0.9848        |
|                | WQR     | 0.3662 | 0.3838        | 0.3832        | 0.3753        | 0.3729        | 0.3718        | 0.3718        | 0.3693        | 0.3693        | <b>0.3862</b>   | <u>0.3845</u> | 0.3810        | 0.3810        |
|                | ACT     | 0.8808 | 0.8897        | 0.8897        | 0.8832        | 0.8841        | 0.8827        | 0.8855        | 0.8827        | 0.8827        | <b>0.8925</b>   | <u>0.8921</u> | 0.8893        | 0.8864        |
|                | CD      | 0.8293 | 0.8285        | 0.8291        | 0.8287        | 0.8285        | 0.8291        | 0.8289        | 0.8294        | 0.8287        | <u>0.8295</u>   | 0.8294        | <b>0.8295</b> | 0.8276        |
| GC             | 0.7450  | 0.7550 | 0.7500        | 0.7650        | 0.7570        | <b>0.7690</b> | 0.7620        | 0.7660        | 0.7630        | 0.7640        | 0.7620          | <u>0.7680</u> | 0.7680        |               |
| Mean           |         | 0.7252 | 0.7400        | <u>0.7411</u> | 0.7380        | 0.7376        | 0.7387        | 0.7378        | 0.7357        | 0.7352        | 0.7408          | <b>0.7412</b> | 0.7392        | 0.7393        |
| Mean Rank      |         | 11.57  | 7.29          | 5.14          | 7.07          | 7.64          | 5.71          | 6.93          | 8.43          | 9.79          | <b>4.14</b>     | <u>4.29</u>   | 6.00          | 7.00          |
| Light-GBM      | AF      | 0.8375 | 0.8285        | 0.8411        | 0.8188        | 0.8244        | 0.8364        | 0.8348        | <b>0.8430</b> | <u>0.8426</u> | 0.8311          | 0.8392        | 0.8366        | 0.8395        |
|                | BH      | 0.5537 | 0.5607        | 0.5636        | <b>0.5693</b> | 0.5618        | 0.5540        | 0.5571        | 0.5478        | 0.5501        | 0.5619          | <u>0.5644</u> | 0.5642        | 0.5595        |
|                | BS      | 0.9429 | 0.9763        | 0.9786        | 0.9751        | 0.9797        | 0.9555        | 0.9565        | 0.9449        | 0.9487        | 0.9737          | 0.9754        | <u>0.9801</u> | <b>0.9813</b> |
|                | WQR     | 0.3825 | 0.4145        | <b>0.4182</b> | 0.3898        | 0.3884        | 0.4131        | 0.4035        | 0.3902        | 0.3952        | 0.4118          | <u>0.4171</u> | 0.4021        | 0.4042        |
|                | ACT     | 0.8832 | 0.8794        | 0.8827        | 0.8808        | 0.8799        | 0.8822        | 0.8860        | 0.8827        | 0.8818        | 0.8888          | <b>0.8925</b> | <u>0.8902</u> | <b>0.8925</b> |
|                | CD      | 0.8300 | 0.8283        | 0.8277        | 0.8293        | 0.8287        | 0.8296        | 0.8298        | 0.8301        | 0.8294        | <u>0.8301</u>   | 0.8297        | <b>0.8303</b> | 0.8294        |
| GC             | 0.7250  | 0.7650 | 0.7600        | 0.7550        | 0.7700        | 0.7490        | 0.7550        | 0.7450        | <u>0.7720</u> | 0.7680        | 0.7720          | <b>0.7760</b> | 0.7700        |               |
| Mean           |         | 0.7364 | 0.7504        | 0.7531        | 0.7454        | 0.7476        | 0.7457        | 0.7461        | 0.7405        | 0.7457        | 0.7522          | <b>0.7558</b> | <u>0.7542</u> | 0.7538        |
| Mean Rank      |         | 9.43   | 8.29          | 5.86          | 8.86          | 8.57          | 8.43          | 7.71          | 8.29          | 7.93          | 5.71            | <b>3.57</b>   | <b>3.57</b>   | <u>4.29</u>   |
| Mean           |         | 0.6806 | 0.7091        | 0.7140        | 0.6953        | 0.6958        | 0.6979        | 0.6976        | 0.6950        | 0.6967        | 0.7171          | <u>0.7185</u> | <b>0.7187</b> | 0.7183        |
| Mean Rank      |         | 10.95  | 7.90          | 5.50          | 8.35          | 8.23          | 7.55          | 7.58          | 8.50          | 8.65          | 4.88            | <b>3.70</b>   | <u>4.25</u>   | 4.98          |

**Table 9:** Performance comparison of PromptFE with and without semantic blinding. For each compared version, the left and middle columns show the results without and with parameter tuning of the downstream model algorithm post AutoFE, respectively, and the right column shows the number of LLM responses. The results where the full version outperforms the blinded version are highlighted in boldface.

| Model          | Dataset | Raw    | GPT-3.5 |        |               |               | GPT-4         |              |        |        |               |               |               |              |
|----------------|---------|--------|---------|--------|---------------|---------------|---------------|--------------|--------|--------|---------------|---------------|---------------|--------------|
|                |         |        | Blinded |        | Full          |               | Blinded       |              | Full   |        |               |               |               |              |
| Linear Model   | AF      | 0.3474 | 0.6613  | 0.6602 | 450.0         | 0.6612        | <b>0.6616</b> | <b>339.8</b> | 0.6678 | 0.6672 | 275.0         | 0.6649        | 0.6647        | 371.4        |
|                | BH      | 0.3776 | 0.4678  | 0.4794 | 438.0         | <b>0.4995</b> | <b>0.5025</b> | <b>378.6</b> | 0.4869 | 0.4996 | 295.6         | <b>0.5184</b> | <b>0.5289</b> | 335.4        |
|                | WQR     | 0.2696 | 0.2643  | 0.2733 | 442.8         | <b>0.2722</b> | <b>0.2745</b> | <b>328.4</b> | 0.2645 | 0.2702 | 244.6         | <b>0.2713</b> | <b>0.2748</b> | 312.6        |
|                | ACT     | 0.8505 | 0.8790  | 0.8799 | 442.8         | 0.8729        | 0.8794        | <b>372.2</b> | 0.8720 | 0.8729 | 238.8         | <b>0.8766</b> | <b>0.8762</b> | 377.4        |
|                | CD      | 0.8267 | 0.8283  | 0.8283 | 454.8         | 0.8282        | 0.8282        | <b>342.0</b> | 0.8282 | 0.8289 | 238.2         | <b>0.8288</b> | 0.8288        | 250.4        |
|                | GC      | 0.7100 | 0.7460  | 0.7390 | 432.2         | <b>0.7570</b> | <b>0.7460</b> | <b>379.0</b> | 0.7430 | 0.7410 | 231.2         | <b>0.7590</b> | <b>0.7420</b> | 310.6        |
| Mean           |         | 0.5636 | 0.6411  | 0.6433 | 443.4         | <b>0.6485</b> | <b>0.6487</b> | <b>356.7</b> | 0.6437 | 0.6461 | 253.9         | <b>0.6532</b> | <b>0.6526</b> | 326.3        |
| Random Forests | AF      | 0.7677 | 0.7644  | 0.7743 | 425.2         | <b>0.7709</b> | <b>0.7787</b> | <b>393.2</b> | 0.7610 | 0.7690 | 274.2         | <b>0.7681</b> | <b>0.7749</b> | 314.2        |
|                | BH      | 0.5426 | 0.5483  | 0.5483 | 479.2         | <b>0.5549</b> | <b>0.5533</b> | <b>374.4</b> | 0.5507 | 0.5491 | 238.4         | <b>0.5543</b> | <b>0.5522</b> | 278.6        |
|                | BS      | 0.9446 | 0.9628  | 0.9628 | 510.0         | <b>0.9873</b> | <b>0.9881</b> | <b>386.8</b> | 0.9535 | 0.9543 | 247.4         | <b>0.9845</b> | <b>0.9848</b> | 255.0        |
|                | WQR     | 0.3662 | 0.3749  | 0.3738 | 461.4         | <b>0.3862</b> | <b>0.3845</b> | <b>362.6</b> | 0.3666 | 0.3674 | 253.0         | <b>0.3810</b> | <b>0.3810</b> | 283.2        |
|                | ACT     | 0.8808 | 0.8864  | 0.8841 | 475.8         | <b>0.8925</b> | <b>0.8921</b> | <b>357.6</b> | 0.8874 | 0.8841 | 222.4         | <b>0.8893</b> | <b>0.8864</b> | 424.0        |
|                | CD      | 0.8293 | 0.8283  | 0.8282 | 497.0         | <b>0.8295</b> | <b>0.8294</b> | <b>349.8</b> | 0.8291 | 0.8286 | 375.2         | <b>0.8295</b> | 0.8276        | <b>304.0</b> |
| GC             | 0.7450  | 0.7630 | 0.7580  | 459.2  | <b>0.7640</b> | <b>0.7620</b> | <b>368.2</b>  | 0.7510       | 0.7490 | 229.6  | <b>0.7680</b> | <b>0.7680</b> | 471.8         |              |
| Mean           |         | 0.6806 | 0.7326  | 0.7328 | 472.5         | <b>0.7408</b> | <b>0.7412</b> | <b>370.4</b> | 0.7285 | 0.7288 | 262.9         | <b>0.7392</b> | <b>0.7393</b> | 333.0        |
| Light-GBM      | AF      | 0.8375 | 0.8304  | 0.8356 | 479.6         | <b>0.8311</b> | <b>0.8392</b> | <b>380.2</b> | 0.8185 | 0.8266 | 284.6         | <b>0.8366</b> | <b>0.8395</b> | 360.6        |
|                | BH      | 0.5537 | 0.5503  | 0.5467 | 490.8         | <b>0.5619</b> | <b>0.5644</b> | <b>342.0</b> | 0.5500 | 0.5609 | 238.4         | <b>0.5642</b> | 0.5595        | 345.6        |
|                | BS      | 0.9429 | 0.9693  | 0.9691 | 480.2         | <b>0.9737</b> | <b>0.9754</b> | <b>380.0</b> | 0.9539 | 0.9536 | 312.6         | <b>0.9801</b> | <b>0.9813</b> | <b>236.8</b> |
|                | WQR     | 0.3825 | 0.4087  | 0.4151 | 493.0         | <b>0.4118</b> | <b>0.4171</b> | <b>322.8</b> | 0.4057 | 0.4050 | 246.8         | 0.4021        | 0.4042        | 293.6        |
|                | ACT     | 0.8832 | 0.8864  | 0.8883 | 513.0         | <b>0.8888</b> | <b>0.8925</b> | <b>367.4</b> | 0.8813 | 0.8748 | 229.0         | <b>0.8902</b> | <b>0.8925</b> | 359.6        |
|                | CD      | 0.8300 | 0.8284  | 0.8292 | 490.8         | <b>0.8301</b> | <b>0.8297</b> | <b>352.2</b> | 0.8295 | 0.8299 | 218.6         | <b>0.8303</b> | 0.8294        | 371.2        |
| GC             | 0.7250  | 0.7620 | 0.7620  | 482.4  | <b>0.7680</b> | <b>0.7720</b> | <b>376.6</b>  | 0.7550       | 0.7550 | 225.0  | <b>0.7760</b> | <b>0.7700</b> | 382.2         |              |
| Mean           |         | 0.6806 | 0.7479  | 0.7494 | 490.0         | <b>0.7522</b> | <b>0.7558</b> | <b>360.2</b> | 0.7420 | 0.7437 | 250.7         | <b>0.7542</b> | <b>0.7538</b> | 335.7        |
| Mean           |         | 0.6806 | 0.7105  | 0.7118 | 469.9         | <b>0.7171</b> | <b>0.7185</b> | <b>362.7</b> | 0.7078 | 0.7092 | 255.9         | <b>0.7187</b> | <b>0.7183</b> | 331.9        |

## D.8 ADDITIONAL HYPERPARAMETER EFFECT

**Number of Iterations.** Figure 14 shows the validation scores on the *AF* and *CD* datasets, which contain the smallest and largest numbers of features, respectively, using Random Forests and Light-GBM. The validation score is evaluated after adding the selected set of candidate features to the dataset, as denoted by  $s_{n^*}$  in line 17 of Algorithm 1. We terminate our algorithm once we have

**Table 10:** Performance comparison of PromptFE with and without RPN canonicalization. For each compared version, the left and middle columns show the results without and with parameter tuning of the downstream model algorithm post AutoFE, respectively, and the right column shows the number of LLM responses. The results where the full version outperforms the reduced version are highlighted in boldface.

| Model          | Dataset | Raw    | GPT-3.5              |        |               |               |               |              | GPT-4                |        |               |               |               |              |
|----------------|---------|--------|----------------------|--------|---------------|---------------|---------------|--------------|----------------------|--------|---------------|---------------|---------------|--------------|
|                |         |        | w/o Canonicalization |        |               | Full          |               |              | w/o Canonicalization |        | Full          |               |               |              |
| Linear Model   | AF      | 0.3474 | 0.6679               | 0.6688 | 338.6         | 0.6612        | 0.6616        | 339.8        | 0.6538               | 0.6529 | 321.2         | <b>0.6649</b> | <b>0.6647</b> | 371.4        |
|                | BH      | 0.3776 | 0.5048               | 0.5076 | 351.2         | 0.4995        | 0.5025        | 378.6        | 0.4987               | 0.5030 | 310.8         | <b>0.5184</b> | <b>0.5289</b> | 335.4        |
|                | WQR     | 0.2696 | 0.2702               | 0.2735 | 336.2         | <b>0.2722</b> | <b>0.2745</b> | <b>328.4</b> | 0.2690               | 0.2706 | 279.0         | <b>0.2713</b> | <b>0.2748</b> | 312.6        |
|                | ACT     | 0.8505 | 0.8748               | 0.8794 | 366.4         | 0.8729        | 0.8794        | 372.2        | 0.8738               | 0.8752 | 298.0         | <b>0.8766</b> | <b>0.8762</b> | 377.4        |
|                | CD      | 0.8267 | 0.8280               | 0.8290 | 350.4         | <b>0.8282</b> | 0.8282        | <b>342.0</b> | 0.8270               | 0.8271 | 285.4         | <b>0.8288</b> | <b>0.8288</b> | <b>250.4</b> |
|                | GC      | 0.7100 | 0.7370               | 0.7330 | 352.0         | <b>0.7570</b> | <b>0.7460</b> | 379.0        | 0.7550               | 0.7490 | 447.2         | <b>0.7590</b> | 0.7420        | <b>310.6</b> |
| Mean           |         | 0.5636 | 0.6471               | 0.6486 | 349.1         | <b>0.6485</b> | <b>0.6487</b> | 356.7        | 0.6462               | 0.6463 | 323.6         | <b>0.6532</b> | <b>0.6526</b> | 326.3        |
| Random Forests | AF      | 0.7677 | 0.7628               | 0.7762 | 358.0         | <b>0.7709</b> | <b>0.7787</b> | 393.2        | 0.7743               | 0.7843 | 340.2         | 0.7681        | 0.7749        | <b>314.2</b> |
|                | BH      | 0.5426 | 0.5573               | 0.5573 | 364.0         | 0.5549        | 0.5533        | 374.4        | 0.5491               | 0.5460 | 322.4         | <b>0.5543</b> | <b>0.5522</b> | <b>278.6</b> |
|                | BS      | 0.9446 | 0.9804               | 0.9807 | 372.2         | <b>0.9873</b> | <b>0.9881</b> | 386.8        | 0.9778               | 0.9777 | 284.4         | <b>0.9845</b> | <b>0.9848</b> | <b>255.0</b> |
|                | WQR     | 0.3662 | 0.3776               | 0.3726 | 334.6         | <b>0.3862</b> | <b>0.3845</b> | 362.6        | 0.3739               | 0.3719 | 269.8         | <b>0.3810</b> | <b>0.3810</b> | 283.2        |
|                | ACT     | 0.8808 | 0.8879               | 0.8841 | 353.4         | <b>0.8925</b> | <b>0.8921</b> | 357.6        | 0.8841               | 0.8864 | 327.6         | <b>0.8893</b> | 0.8864        | 424.0        |
|                | CD      | 0.8293 | 0.8283               | 0.8285 | 381.6         | <b>0.8295</b> | <b>0.8294</b> | <b>349.8</b> | 0.8290               | 0.8287 | 297.2         | <b>0.8295</b> | 0.8276        | 304.0        |
| GC             | 0.7450  | 0.7660 | 0.7620               | 342.2  | 0.7640        | 0.7620        | 368.2         | 0.7680       | 0.7610               | 368.2  | 0.7680        | <b>0.7680</b> | 471.8         |              |
| Mean           |         | 0.6806 | 0.7372               | 0.7373 | 358.0         | <b>0.7408</b> | <b>0.7412</b> | 370.4        | 0.7366               | 0.7366 | 315.7         | <b>0.7392</b> | <b>0.7393</b> | 333.0        |
| Light-GBM      | AF      | 0.8375 | 0.8322               | 0.8365 | 343.6         | 0.8311        | <b>0.8392</b> | 380.2        | 0.8280               | 0.8350 | 376.0         | <b>0.8366</b> | <b>0.8395</b> | <b>360.6</b> |
|                | BH      | 0.5537 | 0.5599               | 0.5556 | 339.2         | <b>0.5619</b> | <b>0.5644</b> | 342.0        | 0.5577               | 0.5548 | 315.2         | <b>0.5642</b> | <b>0.5595</b> | <b>345.6</b> |
|                | BS      | 0.9429 | 0.9643               | 0.9664 | 368.8         | <b>0.9737</b> | <b>0.9754</b> | 380.0        | 0.9597               | 0.9609 | 276.2         | <b>0.9801</b> | <b>0.9813</b> | <b>236.8</b> |
|                | WQR     | 0.3825 | 0.4075               | 0.4042 | 346.4         | <b>0.4118</b> | <b>0.4171</b> | <b>322.8</b> | 0.4036               | 0.4032 | 288.2         | 0.4021        | <b>0.4042</b> | 293.6        |
|                | ACT     | 0.8832 | 0.8813               | 0.8860 | 342.4         | <b>0.8888</b> | <b>0.8925</b> | 367.4        | 0.8822               | 0.8879 | 313.2         | <b>0.8902</b> | <b>0.8925</b> | 359.6        |
|                | CD      | 0.8300 | 0.8302               | 0.8291 | 355.8         | 0.8301        | <b>0.8297</b> | <b>352.2</b> | 0.8295               | 0.8291 | 301.6         | <b>0.8303</b> | <b>0.8294</b> | 371.2        |
| GC             | 0.7250  | 0.7640 | 0.7650               | 346.2  | <b>0.7680</b> | <b>0.7720</b> | 376.6         | 0.7620       | 0.7650               | 428.8  | <b>0.7760</b> | <b>0.7700</b> | <b>382.2</b>  |              |
| Mean           |         | 0.6806 | 0.7485               | 0.7490 | 348.9         | <b>0.7522</b> | <b>0.7558</b> | 360.2        | 0.7461               | 0.7480 | 328.5         | <b>0.7542</b> | <b>0.7538</b> | 335.7        |
| Mean           |         | 0.6806 | 0.7141               | 0.7148 | 352.2         | <b>0.7171</b> | <b>0.7185</b> | 362.7        | 0.7128               | 0.7135 | 322.5         | <b>0.7187</b> | <b>0.7183</b> | 331.9        |

**Table 11:** Standard deviations of Table 1, summary of experimental results.

| Model          | Dataset | Raw    | DIFER  |        | OpenFE |        | CAAFE   |        |        |        | FEBP (ours) |        |        |        |
|----------------|---------|--------|--------|--------|--------|--------|---------|--------|--------|--------|-------------|--------|--------|--------|
|                |         |        |        |        |        |        | GPT-3.5 |        | GPT-4  |        | GPT-3.5     |        | GPT-4  |        |
| Linear Model   | AF      | -      | 0.2559 | 0.2012 | 0.0015 | 0.0014 | 0.0099  | 0.0102 | 0.0511 | 0.0513 | 0.0101      | 0.0100 | 0.0267 | 0.0268 |
|                | BH      | -      | 0.0092 | 0.0153 | 0.0169 | 0.0188 | 0.0196  | 0.0184 | 0.0408 | 0.0419 | 0.0111      | 0.0149 | 0.0254 | 0.0184 |
|                | WQR     | -      | 0.0305 | 0.0223 | 0.0058 | 0.0055 | 0.0046  | 0.0038 | 0.0060 | 0.0060 | 0.0135      | 0.0112 | 0.0068 | 0.0044 |
|                | ACT     | -      | 0.0179 | 0.0073 | 0.0140 | 0.0105 | 0.0035  | 0.0021 | 0.0054 | 0.0053 | 0.0085      | 0.0051 | 0.0040 | 0.0062 |
|                | CD      | -      | 0.0014 | 0.0006 | 0.0006 | 0.0002 | 0.0006  | 0.0007 | 0.0057 | 0.0051 | 0.0013      | 0.0007 | 0.0006 | 0.0009 |
|                | GC      | -      | 0.0272 | 0.0104 | 0.0097 | 0.0076 | 0.0100  | 0.0125 | 0.0134 | 0.0108 | 0.0120      | 0.0213 | 0.0108 | 0.0152 |
| Random Forests | AF      | -      | 0.0054 | 0.0044 | 0.0032 | 0.0036 | 0.0032  | 0.0034 | 0.0108 | 0.0084 | 0.0090      | 0.0086 | 0.0059 | 0.0095 |
|                | BH      | -      | 0.0142 | 0.0131 | 0.0034 | 0.0068 | 0.0050  | 0.0050 | 0.0084 | 0.0113 | 0.0057      | 0.0077 | 0.0059 | 0.0046 |
|                | BS      | -      | 0.0128 | 0.0113 | 0.0003 | 0.0003 | 0.0003  | 0.0003 | 0.0208 | 0.0207 | 0.0088      | 0.0070 | 0.0157 | 0.0154 |
|                | WQR     | -      | 0.0108 | 0.0109 | 0.0030 | 0.0076 | 0.0022  | 0.0022 | 0.0051 | 0.0051 | 0.0034      | 0.0069 | 0.0022 | 0.0026 |
|                | ACT     | -      | 0.0048 | 0.0058 | 0.0037 | 0.0087 | 0.0030  | 0.0055 | 0.0020 | 0.0030 | 0.0055      | 0.0051 | 0.0043 | 0.0054 |
|                | CD      | -      | 0.0010 | 0.0011 | 0.0003 | 0.0004 | 0.0005  | 0.0004 | 0.0008 | 0.0001 | 0.0011      | 0.0010 | 0.0009 | 0.0017 |
| GC             | -       | 0.0184 | 0.0177 | 0.0154 | 0.0110 | 0.0082 | 0.0076  | 0.0065 | 0.0164 | 0.0114 | 0.0067      | 0.0097 | 0.0097 |        |
| Light-GBM      | AF      | -      | 0.0029 | 0.0029 | 0.0058 | 0.0036 | 0.0067  | 0.0027 | 0.0072 | 0.0077 | 0.0129      | 0.0054 | 0.0061 | 0.0041 |
|                | BH      | -      | 0.0147 | 0.0260 | 0.0128 | 0.0150 | 0.0114  | 0.0111 | 0.0145 | 0.0188 | 0.0169      | 0.0076 | 0.0134 | 0.0073 |
|                | BS      | -      | 0.0092 | 0.0070 | 0.0007 | 0.0004 | 0.0159  | 0.0198 | 0.0056 | 0.0139 | 0.0151      | 0.0139 | 0.0033 | 0.0034 |
|                | WQR     | -      | 0.0134 | 0.0164 | 0.0072 | 0.0133 | 0.0084  | 0.0080 | 0.0116 | 0.0134 | 0.0123      | 0.0085 | 0.0097 | 0.0092 |
|                | ACT     | -      | 0.0048 | 0.0042 | 0.0068 | 0.0094 | 0.0061  | 0.0045 | 0.0045 | 0.0027 | 0.0027      | 0.0017 | 0.0050 | 0.0077 |
|                | CD      | -      | 0.0009 | 0.0013 | 0.0004 | 0.0010 | 0.0008  | 0.0005 | 0.0010 | 0.0007 | 0.0004      | 0.0004 | 0.0004 | 0.0008 |
| GC             | -       | 0.0141 | 0.0184 | 0.0184 | 0.0184 | 0.0222 | 0.0166  | 0.0079 | 0.0199 | 0.0076 | 0.0045      | 0.0096 | 0.0146 |        |

200 candidate features, as constructing additional features does not substantially enhance the performance, but constructing fewer features degrades the performance in some cases.

**Number of Examples in Prompt with GPT-4.** Table 21 reports the maximum validation score across iterations by varying the number of example features provided in the prompt to GPT-4. We observe improved performance as the number of example features increases. This suggests that providing more example features helps fully leverage GPT-4’s enhanced in-context learning capabilities. In our experiments, we set the number of example features to 10 for a fair comparison with GPT-3.5.

**Table 12:** Standard deviations of Table 2, performance comparison of FEBP with and without semantic blinding.

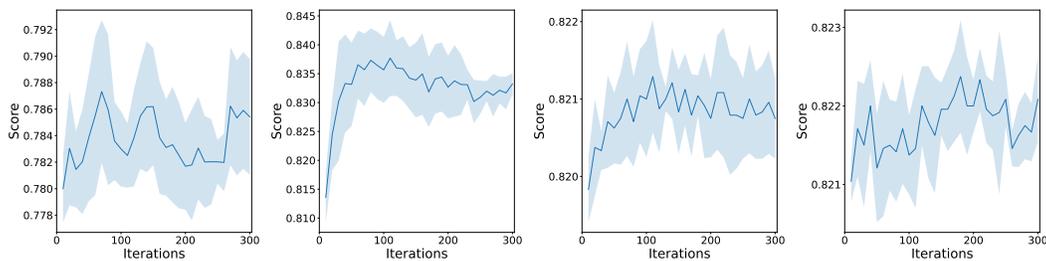
| Model          | Dataset | Raw    | GPT-3.5 |        |        |        |        |        | GPT-4   |        |        |        |        |       |
|----------------|---------|--------|---------|--------|--------|--------|--------|--------|---------|--------|--------|--------|--------|-------|
|                |         |        | Blinded |        |        | Full   |        |        | Blinded |        |        | Full   |        |       |
| Linear Model   | AF      | -      | 0.0147  | 0.0156 | 36.1   | 0.0101 | 0.0100 | 28.8   | 0.0162  | 0.0161 | 25.8   | 0.0267 | 0.0268 | 92.3  |
|                | BH      | -      | 0.0444  | 0.0519 | 39.0   | 0.0111 | 0.0149 | 42.2   | 0.0161  | 0.0131 | 66.7   | 0.0254 | 0.0184 | 58.6  |
|                | WQR     | -      | 0.0133  | 0.0032 | 48.9   | 0.0135 | 0.0112 | 15.3   | 0.0128  | 0.0046 | 23.5   | 0.0068 | 0.0044 | 80.6  |
|                | ACT     | -      | 0.0088  | 0.0107 | 15.4   | 0.0085 | 0.0051 | 17.5   | 0.0056  | 0.0085 | 15.5   | 0.0040 | 0.0062 | 54.8  |
|                | CD      | -      | 0.0014  | 0.0003 | 27.6   | 0.0013 | 0.0007 | 13.1   | 0.0021  | 0.0011 | 13.2   | 0.0006 | 0.0009 | 14.8  |
|                | GC      | -      | 0.0114  | 0.0042 | 32.3   | 0.0120 | 0.0213 | 14.3   | 0.0125  | 0.0114 | 11.0   | 0.0108 | 0.0152 | 36.4  |
| Random Forests | AF      | -      | 0.0086  | 0.0058 | 60.3   | 0.0090 | 0.0086 | 47.3   | 0.0092  | 0.0079 | 27.9   | 0.0059 | 0.0095 | 93.6  |
|                | BH      | -      | 0.0068  | 0.0068 | 45.3   | 0.0057 | 0.0077 | 14.5   | 0.0142  | 0.0132 | 24.7   | 0.0059 | 0.0046 | 23.0  |
|                | BS      | -      | 0.0186  | 0.0181 | 112.1  | 0.0088 | 0.0070 | 47.8   | 0.0103  | 0.0088 | 38.8   | 0.0157 | 0.0154 | 39.2  |
|                | WQR     | -      | 0.0078  | 0.0081 | 40.5   | 0.0034 | 0.0069 | 18.5   | 0.0092  | 0.0075 | 19.1   | 0.0022 | 0.0026 | 45.2  |
|                | ACT     | -      | 0.0099  | 0.0035 | 33.7   | 0.0055 | 0.0051 | 13.1   | 0.0100  | 0.0093 | 16.6   | 0.0043 | 0.0054 | 85.7  |
|                | CD      | -      | 0.0015  | 0.0008 | 53.3   | 0.0011 | 0.0010 | 14.5   | 0.0005  | 0.0008 | 83.4   | 0.0009 | 0.0017 | 56.9  |
| Light-GBM      | GC      | -      | 0.0067  | 0.0057 | 28.9   | 0.0114 | 0.0067 | 17.3   | 0.0210  | 0.0143 | 12.8   | 0.0097 | 0.0097 | 113.1 |
|                | AF      | -      | 0.0104  | 0.0060 | 66.8   | 0.0129 | 0.0054 | 21.7   | 0.0142  | 0.0155 | 39.6   | 0.0061 | 0.0041 | 73.1  |
|                | BH      | -      | 0.0131  | 0.0170 | 60.7   | 0.0169 | 0.0076 | 20.7   | 0.0119  | 0.0121 | 25.7   | 0.0134 | 0.0073 | 36.1  |
|                | BS      | -      | 0.0152  | 0.0178 | 76.3   | 0.0151 | 0.0139 | 31.8   | 0.0048  | 0.0049 | 74.5   | 0.0033 | 0.0034 | 32.1  |
|                | WQR     | -      | 0.0151  | 0.0028 | 36.9   | 0.0123 | 0.0085 | 17.3   | 0.0195  | 0.0190 | 21.1   | 0.0097 | 0.0092 | 46.3  |
|                | ACT     | -      | 0.0021  | 0.0030 | 44.2   | 0.0027 | 0.0017 | 28.5   | 0.0042  | 0.0128 | 15.7   | 0.0050 | 0.0077 | 49.6  |
| CD             | -       | 0.0011 | 0.0011  | 59.4   | 0.0004 | 0.0004 | 15.7   | 0.0007 | 0.0010  | 5.6    | 0.0004 | 0.0008 | 85.7   |       |
|                | GC      | -      | 0.0130  | 0.0148 | 41.7   | 0.0076 | 0.0045 | 23.0   | 0.0117  | 0.0094 | 13.7   | 0.0096 | 0.0146 | 46.9  |

**Table 13:** Standard deviations of Table 4, effect of temperature.

| Model | Dataset | Temperature |       |        |      |        |      |
|-------|---------|-------------|-------|--------|------|--------|------|
|       |         | 0.5         |       | 1      |      | 1.5    |      |
| RF    | AF      | 0.0071      | 160.9 | 0.0042 | 47.3 | 0.0040 | 34.7 |
|       | CD      | 0.0005      | 324.3 | 0.0004 | 14.5 | 0.0005 | 64.1 |
| LGBM  | AF      | 0.0042      | 523.3 | 0.0044 | 21.7 | 0.0022 | 59.8 |
|       | CD      | 0.0008      | 174.7 | 0.0007 | 15.7 | 0.0005 | 73.0 |

**Table 14:** Standard deviations of Table 3, effect of the number of example features in the prompt.

| Model | Dataset | Number of Examples |       |        |      |        |      |        |      |
|-------|---------|--------------------|-------|--------|------|--------|------|--------|------|
|       |         | 1                  |       | 5      |      | 10     |      | 20     |      |
| RF    | AF      | 0.0054             | 55.8  | 0.0035 | 45.0 | 0.0042 | 47.3 | 0.0056 | 24.0 |
|       | WQR     | 0.0088             | 19.6  | 0.0038 | 11.4 | 0.0027 | 18.5 | 0.0096 | 29.6 |
|       | CD      | 0.0005             | 46.5  | 0.0007 | 19.1 | 0.0004 | 14.5 | 0.0006 | 17.8 |
| LGBM  | AF      | 0.0065             | 103.2 | 0.0031 | 21.6 | 0.0044 | 21.7 | 0.0044 | 56.4 |
|       | WQR     | 0.0048             | 16.9  | 0.0057 | 32.4 | 0.0064 | 17.3 | 0.0064 | 26.5 |
|       | CD      | 0.0003             | 71.2  | 0.0002 | 39.0 | 0.0007 | 15.7 | 0.0005 | 17.5 |

**Figure 14:** The validation score across iterations using Random Forests and LightGBM.

**Table 15:** The percentage performance improvement of FEBP over the baseline methods, with GPT-3.5. For each compared method, the left and right columns show the results without and with parameter tuning of the downstream model algorithm post AutoFE, respectively.

| Model          | Dataset | Raw   |       | DIFER |       | OpenFE |       | CAAFE |       |
|----------------|---------|-------|-------|-------|-------|--------|-------|-------|-------|
| Linear Model   | AF      | 90.34 | 90.46 | 12.65 | 8.64  | 53.77  | 53.77 | 64.86 | 64.76 |
|                | BH      | 32.27 | 33.06 | -0.37 | 0.61  | 28.06  | 29.51 | 4.32  | 5.46  |
|                | WQR     | 0.96  | 1.80  | 9.97  | 4.37  | 0.35   | 0.32  | -0.74 | -0.46 |
|                | ACT     | 2.64  | 3.41  | 0.16  | -0.05 | 0.00   | 0.75  | 2.47  | 3.29  |
|                | CD      | 0.18  | 0.19  | 0.10  | 0.03  | 0.21   | 0.17  | 0.20  | 0.19  |
|                | GC      | 6.62  | 5.07  | 6.02  | 0.54  | 3.42   | 2.47  | 2.99  | 1.77  |
| Mean           |         | 22.17 | 22.33 | 4.76  | 2.36  | 14.30  | 14.50 | 12.35 | 12.50 |
| Random Forests | AF      | 0.42  | 1.44  | 0.78  | 0.02  | 1.72   | 1.37  | -0.02 | 1.23  |
|                | BH      | 2.26  | 1.97  | -2.95 | -2.95 | -1.92  | -1.55 | -0.13 | -0.41 |
|                | BS      | 4.52  | 4.60  | 0.08  | 0.10  | -0.29  | -0.21 | -0.43 | -0.35 |
|                | WQR     | 5.44  | 5.00  | 0.61  | 0.34  | 2.89   | 3.11  | 3.86  | 3.42  |
|                | ACT     | 1.33  | 1.27  | 0.32  | 0.26  | 1.06   | 0.90  | 1.11  | 0.74  |
|                | CD      | 0.02  | 0.01  | 0.12  | 0.04  | 0.09   | 0.10  | 0.05  | 0.06  |
| Mean           |         | 2.55  | 2.28  | 1.19  | 1.60  | -0.13  | 0.66  | -0.65 | 0.00  |
| Light-GBM      | AF      | -0.76 | 0.20  | 0.32  | -0.23 | 1.51   | 1.80  | -0.63 | 0.53  |
|                | BH      | 1.48  | 1.94  | 0.21  | 0.14  | -1.30  | 0.47  | 1.42  | 1.32  |
|                | BS      | 3.27  | 3.45  | -0.27 | -0.32 | -0.15  | -0.43 | 1.91  | 1.98  |
|                | WQR     | 7.67  | 9.04  | -0.63 | -0.27 | 5.66   | 7.40  | -0.29 | 3.36  |
|                | ACT     | 0.63  | 1.06  | 1.06  | 1.11  | 0.90   | 1.43  | 0.74  | 0.74  |
|                | CD      | 0.02  | -0.04 | 0.22  | 0.25  | 0.10   | 0.12  | 0.06  | -0.01 |
| Mean           |         | 5.93  | 6.48  | 0.39  | 1.58  | 1.72   | 0.26  | 2.54  | 2.25  |
| Mean           |         | 2.60  | 3.16  | 0.18  | 0.32  | 1.21   | 1.58  | 0.82  | 1.45  |
| Mean           |         | 8.39  | 8.63  | 1.50  | 0.79  | 4.88   | 5.12  | 4.18  | 4.49  |

**Table 16:** The percentage performance improvement of FEBP over the baseline methods, with GPT-4. For each compared method, the left and right columns show the results without and with parameter tuning of the downstream model algorithm post AutoFE, respectively.

| Model          | Dataset | Raw   |       | DIFER |       | OpenFE |       | CAAFE |       |
|----------------|---------|-------|-------|-------|-------|--------|-------|-------|-------|
| Linear Model   | AF      | 91.40 | 91.34 | 13.27 | 9.14  | 54.62  | 54.48 | 51.94 | 51.82 |
|                | BH      | 37.28 | 40.06 | 3.41  | 5.90  | 32.92  | 36.33 | 15.14 | 17.39 |
|                | WQR     | 0.64  | 1.94  | 9.62  | 4.51  | 0.03   | 0.46  | -2.25 | -0.99 |
|                | ACT     | 3.08  | 3.02  | 0.59  | -0.42 | 0.43   | 0.37  | 2.35  | 2.24  |
|                | CD      | 0.26  | 0.26  | 0.18  | 0.10  | 0.28   | 0.25  | 0.61  | 0.61  |
|                | GC      | 6.90  | 4.51  | 6.30  | 0.00  | 3.69   | 1.92  | 5.27  | 2.91  |
| Mean           |         | 23.26 | 23.52 | 5.56  | 3.21  | 15.33  | 15.64 | 12.17 | 12.33 |
| Random Forests | AF      | 0.05  | 0.93  | 0.41  | -0.47 | 1.35   | 0.86  | -0.20 | 0.37  |
|                | BH      | 2.16  | 1.76  | -3.05 | -3.15 | -2.02  | -1.76 | 0.56  | 0.54  |
|                | BS      | 4.23  | 4.25  | -0.20 | -0.23 | -0.57  | -0.54 | 0.28  | 0.32  |
|                | WQR     | 4.03  | 4.03  | -0.74 | -0.58 | 1.51   | 2.16  | 3.17  | 3.17  |
|                | ACT     | 0.95  | 0.64  | -0.05 | -0.37 | 0.69   | 0.26  | 0.74  | 0.42  |
|                | CD      | 0.02  | -0.20 | 0.13  | -0.18 | 0.10   | -0.11 | 0.02  | -0.13 |
| Mean           |         | 3.09  | 3.09  | 1.72  | 2.40  | 0.39   | 1.45  | 0.26  | 0.66  |
| Mean           |         | 2.08  | 2.07  | -0.26 | -0.37 | 0.21   | 0.33  | 0.69  | 0.76  |
| Light-GBM      | AF      | -0.11 | 0.24  | 0.98  | -0.19 | 2.18   | 1.84  | -0.75 | -0.36 |
|                | BH      | 1.90  | 1.04  | 0.63  | -0.74 | -0.89  | -0.41 | 3.00  | 1.70  |
|                | BS      | 3.94  | 4.08  | 0.38  | 0.28  | 0.51   | 0.17  | 3.72  | 3.44  |
|                | WQR     | 5.12  | 5.67  | -2.98 | -3.35 | 3.16   | 4.08  | 3.04  | 2.28  |
|                | ACT     | 0.79  | 1.06  | 1.22  | 1.11  | 1.06   | 1.43  | 0.85  | 1.22  |
|                | CD      | 0.04  | -0.07 | 0.24  | 0.21  | 0.12   | 0.08  | 0.03  | 0.00  |
| Mean           |         | 7.03  | 6.21  | 1.44  | 1.32  | 2.78   | 0.00  | 4.16  | -0.26 |
| Mean           |         | 2.67  | 2.60  | 0.27  | -0.19 | 1.27   | 1.03  | 2.01  | 1.15  |
| Mean           |         | 8.64  | 8.69  | 1.67  | 0.76  | 5.12   | 5.17  | 4.60  | 4.37  |

**Table 17:** The Nemenyi post-hoc test  $p$ -values for pairwise comparison of the methods in Table 1. Results that are significant at the  $p = 0.05$  confidence level are highlighted in boldface.

|         | Raw           | DIFER         | OpenFE        | CAAFFE        |               | FEBP (ours)   |               |               |               |               |               |               |               |
|---------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|
|         |               |               |               | GPT-3.5       | GPT-4         | GPT-3.5       | GPT-4         |               |               |               |               |               |               |
| Raw     | 1.0000        | <b>0.0010</b> | <b>0.0010</b> | <b>0.0203</b> | <b>0.0086</b> | <b>0.0010</b> | <b>0.0010</b> | <b>0.0409</b> | <b>0.0179</b> | <b>0.0010</b> | <b>0.0010</b> | <b>0.0010</b> | <b>0.0010</b> |
| DIFER   | <b>0.0010</b> | 1.0000        | 0.3235        | 0.4051        | 0.5626        | 0.9000        | 0.9000        | 0.2697        | 0.4310        | <b>0.0397</b> | <b>0.0010</b> | <b>0.0028</b> | 0.1535        |
| OpenFE  | <b>0.0203</b> | 0.4051        | <b>0.0010</b> | 1.0000        | 0.9000        | 0.9000        | 0.9000        | 0.9000        | 0.9000        | <b>0.0010</b> | <b>0.0010</b> | <b>0.0010</b> | <b>0.0010</b> |
| CAAFFE  | <b>0.0086</b> | 0.5626        | <b>0.0010</b> | 0.9000        | 1.0000        | 0.9000        | 0.9000        | 0.9000        | 0.9000        | <b>0.0010</b> | <b>0.0010</b> | <b>0.0010</b> | <b>0.0010</b> |
| GPT-3.5 | <b>0.0010</b> | 0.9000        | <b>0.0343</b> | 0.9000        | 0.9000        | 1.0000        | 0.9000        | 0.8216        | 0.9000        | <b>0.0016</b> | <b>0.0010</b> | <b>0.0010</b> | <b>0.0105</b> |
| GPT-4   | <b>0.0010</b> | 0.9000        | <b>0.0397</b> | 0.9000        | 0.9000        | 0.9000        | 1.0000        | 0.7929        | 0.9000        | <b>0.0019</b> | <b>0.0010</b> | <b>0.0010</b> | <b>0.0125</b> |
| FEBP    | <b>0.0409</b> | 0.2697        | <b>0.0010</b> | 0.9000        | 0.9000        | 0.8216        | 0.7929        | 1.0000        | 0.9000        | <b>0.0010</b> | <b>0.0010</b> | <b>0.0010</b> | <b>0.0010</b> |
| GPT-3.5 | <b>0.0179</b> | 0.4310        | <b>0.0010</b> | 0.9000        | 0.9000        | 0.9000        | 0.9000        | 1.0000        | 1.0000        | <b>0.0010</b> | <b>0.0010</b> | <b>0.0010</b> | <b>0.0010</b> |
| GPT-4   | <b>0.0010</b> | <b>0.0397</b> | 0.9000        | <b>0.0010</b> | <b>0.0010</b> | <b>0.0016</b> | <b>0.0019</b> | <b>0.0010</b> | <b>0.0010</b> | 0.9000        | 0.9000        | 0.9000        | 0.9000        |
| GPT-3.5 | <b>0.0010</b> | <b>0.0010</b> | 0.7526        | <b>0.0010</b> | <b>0.0010</b> | <b>0.0010</b> | <b>0.0010</b> | <b>0.0010</b> | <b>0.0010</b> | 0.9000        | 1.0000        | 0.9000        | 0.9000        |
| GPT-4   | <b>0.0010</b> | <b>0.0028</b> | 0.9000        | <b>0.0010</b> | <b>0.0010</b> | <b>0.0010</b> | <b>0.0010</b> | <b>0.0010</b> | <b>0.0010</b> | 0.9000        | 0.9000        | 1.0000        | 0.9000        |
| GPT-4   | <b>0.0010</b> | 0.1535        | 0.9000        | <b>0.0010</b> | <b>0.0010</b> | <b>0.0105</b> | <b>0.0125</b> | <b>0.0010</b> | <b>0.0010</b> | 0.9000        | 0.9000        | 0.9000        | 1.0000        |

**Table 18:** The Nemenyi post-hoc test  $p$ -values for pairwise comparison of the methods in Table 2. Results that are significant at the  $p = 0.05$  confidence level are highlighted in boldface.

|         | Raw           | GPT-3.5       |               |               | GPT-4         |               |               |               |               |
|---------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|
|         |               | Blinded       | Full          | Blinded       | Full          | Blinded       | Full          |               |               |
| Raw     | 1.0000        | <b>0.0010</b> | <b>0.0010</b> | <b>0.0010</b> | <b>0.0010</b> | <b>0.0017</b> | <b>0.0010</b> | <b>0.0010</b> | <b>0.0010</b> |
| Blinded | <b>0.0010</b> | 1.0000        | 0.9000        | <b>0.0062</b> | <b>0.0010</b> | 0.9000        | 0.9000        | <b>0.0010</b> | <b>0.0057</b> |
| GPT-3.5 | <b>0.0010</b> | 0.9000        | 1.0000        | 0.1775        | <b>0.0066</b> | 0.3858        | 0.9000        | <b>0.0105</b> | 0.1677        |
| Full    | <b>0.0010</b> | <b>0.0062</b> | 0.1775        | 1.0000        | 0.9000        | <b>0.0010</b> | <b>0.0069</b> | 0.9000        | 0.9000        |
| Blinded | <b>0.0017</b> | 0.9000        | 0.3858        | <b>0.0010</b> | <b>0.0010</b> | 1.0000        | 0.9000        | <b>0.0010</b> | <b>0.0010</b> |
| GPT-4   | <b>0.0010</b> | 0.9000        | 0.9000        | <b>0.0069</b> | <b>0.0010</b> | 0.9000        | 1.0000        | <b>0.0010</b> | <b>0.0062</b> |
| Full    | <b>0.0010</b> | <b>0.0010</b> | <b>0.0105</b> | 0.9000        | 0.9000        | <b>0.0010</b> | <b>0.0010</b> | 1.0000        | 0.9000        |
| Full    | <b>0.0010</b> | <b>0.0057</b> | 0.1677        | 0.9000        | 0.9000        | <b>0.0010</b> | <b>0.0062</b> | 0.9000        | 1.0000        |

**Table 19:** The Nemenyi post-hoc test  $p$ -values for pairwise comparison of the methods in Table 1, excluding linear model results. Results that are significant at the  $p = 0.05$  confidence level are highlighted in boldface.

|         | Raw           | DIFER         | OpenFE        | CAAFFE        |               | FEBP (ours)   |               |               |               |               |               |               |               |
|---------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|
|         |               |               |               | GPT-3.5       | GPT-4         | GPT-3.5       | GPT-4         |               |               |               |               |               |               |
| Raw     | 1.0000        | <b>0.0010</b> | <b>0.0010</b> | 0.5006        | 0.3953        | <b>0.0010</b> | <b>0.0012</b> | 0.3875        | 0.2344        | <b>0.0010</b> | <b>0.0010</b> | <b>0.0010</b> | <b>0.0010</b> |
| DIFER   | <b>0.0010</b> | 1.0000        | 0.9000        | 0.6382        | 0.7345        | 0.9000        | 0.9000        | 0.7414        | 0.8996        | 0.4263        | <b>0.0299</b> | 0.1392        | 0.9000        |
| OpenFE  | 0.5006        | 0.6382        | <b>0.0098</b> | 1.0000        | 0.9000        | 0.6175        | 0.7138        | 0.9000        | 0.9000        | <b>0.0010</b> | <b>0.0010</b> | <b>0.0010</b> | <b>0.0171</b> |
| CAAFFE  | 0.3953        | 0.7345        | <b>0.0171</b> | 0.9000        | 1.0000        | 0.7138        | 0.8102        | 0.9000        | 0.9000        | <b>0.0010</b> | <b>0.0010</b> | <b>0.0010</b> | <b>0.0289</b> |
| GPT-3.5 | <b>0.0010</b> | 0.9000        | 0.9000        | 0.6175        | 0.7138        | 1.0000        | 0.9000        | 0.7207        | 0.8790        | 0.4493        | <b>0.0334</b> | 0.1516        | 0.9000        |
| GPT-4   | <b>0.0012</b> | 0.9000        | 0.8308        | 0.7138        | 0.8102        | 0.9000        | 1.0000        | 0.8170        | 0.9000        | 0.3422        | <b>0.0199</b> | 0.1017        | 0.9000        |
| FEBP    | 0.3875        | 0.7414        | <b>0.0178</b> | 0.9000        | 0.9000        | 0.7207        | 0.8170        | 1.0000        | 0.9000        | <b>0.0010</b> | <b>0.0010</b> | <b>0.0010</b> | <b>0.0299</b> |
| GPT-3.5 | 0.2344        | 0.8996        | <b>0.0412</b> | 0.9000        | 0.9000        | 0.8790        | 0.9000        | 0.9000        | 1.0000        | <b>0.0026</b> | <b>0.0010</b> | <b>0.0010</b> | 0.0661        |
| GPT-4   | <b>0.0010</b> | 0.4263        | 0.9000        | <b>0.0010</b> | <b>0.0010</b> | 0.4493        | 0.3422        | <b>0.0010</b> | <b>0.0026</b> | 1.0000        | 0.9000        | 0.9000        | 0.9000        |
| GPT-3.5 | <b>0.0010</b> | <b>0.0299</b> | 0.8377        | <b>0.0010</b> | <b>0.0010</b> | <b>0.0334</b> | <b>0.0199</b> | <b>0.0010</b> | <b>0.0010</b> | 0.9000        | 1.0000        | 0.9000        | 0.7414        |
| GPT-4   | <b>0.0010</b> | 0.1392        | 0.9000        | <b>0.0010</b> | <b>0.0010</b> | 0.1516        | 0.1017        | <b>0.0010</b> | <b>0.0010</b> | 0.9000        | 0.9000        | 1.0000        | 0.9000        |
| GPT-4   | <b>0.0010</b> | 0.9000        | 0.9000        | <b>0.0171</b> | <b>0.0289</b> | 0.9000        | 0.9000        | <b>0.0299</b> | 0.0661        | 0.9000        | 0.7414        | 0.9000        | 1.0000        |

**Table 20:** The Nemenyi post-hoc test  $p$ -values for pairwise comparison of the methods in Table 10. Results that are significant at the  $p = 0.05$  confidence level are highlighted in boldface.

|         | Raw           | GPT-3.5       |               |               | GPT-4         |               |               |               |               |
|---------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|---------------|
|         |               | w/o           | Full          | w/o           | Full          | w/o           | Full          |               |               |
| Raw     | 1.0000        | <b>0.0010</b> |
| w/o     | <b>0.0010</b> | 1.0000        | 0.9000        | 0.2977        | <b>0.0060</b> | 0.9000        | 0.9000        | <b>0.0224</b> | 0.4293        |
| GPT-3.5 | <b>0.0010</b> | 0.9000        | 1.0000        | 0.6618        | <b>0.0433</b> | 0.8811        | 0.8889        | 0.1230        | 0.7871        |
| Full    | <b>0.0010</b> | 0.2977        | 0.6618        | 1.0000        | 0.9000        | <b>0.0341</b> | <b>0.0355</b> | 0.9000        | 0.9000        |
| w/o     | <b>0.0010</b> | <b>0.0060</b> | <b>0.0433</b> | 0.9000        | 1.0000        | <b>0.0010</b> | <b>0.0010</b> | 0.9000        | 0.8028        |
| GPT-4   | <b>0.0010</b> | 0.9000        | 0.8811        | <b>0.0341</b> | <b>0.0010</b> | 1.0000        | 0.9000        | <b>0.0010</b> | 0.0635        |
| w/o     | <b>0.0010</b> | 0.9000        | 0.8889        | <b>0.0355</b> | <b>0.0010</b> | 0.9000        | 1.0000        | <b>0.0010</b> | 0.0659        |
| Full    | <b>0.0010</b> | <b>0.0224</b> | 0.1230        | 0.9000        | 0.9000        | <b>0.0010</b> | <b>0.0010</b> | 1.0000        | 0.9000        |
| Full    | <b>0.0010</b> | 0.4293        | 0.7871        | 0.9000        | 0.8028        | 0.0635        | 0.0659        | 0.9000        | 1.0000        |

1242  
 1243  
 1244  
 1245  
 1246  
 1247  
 1248  
 1249  
 1250  
 1251  
 1252  
 1253  
 1254  
 1255  
 1256  
 1257  
 1258  
 1259  
 1260  
 1261  
 1262  
 1263  
 1264  
 1265  
 1266  
 1267  
 1268  
 1269  
 1270  
 1271  
 1272  
 1273  
 1274  
 1275  
 1276  
 1277  
 1278  
 1279  
 1280  
 1281  
 1282  
 1283  
 1284  
 1285  
 1286  
 1287  
 1288  
 1289  
 1290  
 1291  
 1292  
 1293  
 1294  
 1295

**Table 21:** Effect of the number of example features in the prompt with GPT-4. For each compared setting, the left column shows the validation score, and the right column shows the number of LLM responses. The best results are boldfaced.

| Model | Dataset | Number of Examples |               |               |               |               |
|-------|---------|--------------------|---------------|---------------|---------------|---------------|
|       |         | 1                  | 5             | 10            | 20            | 30            |
| RF    | AF      | 0.7864             | <b>0.7922</b> | 0.7905        | 0.7897        | 0.7920        |
|       | WQR     | 0.3847             | 0.3835        | 0.3839        | 0.3862        | <b>0.3862</b> |
|       | CD      | 0.8219             | 0.8218        | 0.8218        | 0.8219        | <b>0.8222</b> |
| LGBM  | AF      | 0.8387             | 0.8413        | 0.8401        | <b>0.8433</b> | 0.8411        |
|       | WQR     | 0.4216             | 0.4242        | <b>0.4290</b> | 0.4258        | 0.4267        |
|       | CD      | 0.8231             | <b>0.8234</b> | 0.8227        | 0.8229        | 0.8231        |
| Mean  |         | 0.6794             | 0.6810        | 0.6813        | 0.6816        | <b>0.6819</b> |

1296 D.9 NUMBER OF SELECTED FEATURES  
1297

1298 Table 22 compares the number of features added to the datasets. Our method FEBP adaptively de-  
1299 termines the number of features and selects fewer features than DIFER (Zhu et al., 2022b), demon-  
1300 strating the effectiveness of the features generated by our method.

1301 **Table 22:** Comparison of the number of selected features.  
1302

1303

| Model          | Dataset | DIFER | OpenFE | FEBP Blinded |       | FEBP    |       |
|----------------|---------|-------|--------|--------------|-------|---------|-------|
|                |         |       |        | GPT-3.5      | GPT-4 | GPT-3.5 | GPT-4 |
| Linear Model   | AF      | 310   | 10     | 167          | 165   | 162     | 183   |
|                | BH      | 156   | 10     | 104          | 141   | 144     | 90    |
|                | WQR     | 109   | 10     | 57           | 80    | 43      | 55    |
|                | ACT     | 113   | 10     | 84           | 49    | 85      | 14    |
|                | CD      | 157   | 10     | 92           | 68    | 74      | 74    |
|                | GC      | 105   | 10     | 75           | 97    | 120     | 51    |
| Random Forests | AF      | 387   | 10     | 39           | 19    | 15      | 34    |
|                | BH      | 186   | 10     | 4            | 6     | 19      | 77    |
|                | BS      | 46    | 10     | 9            | 7     | 9       | 65    |
|                | WQR     | 63    | 10     | 9            | 44    | 39      | 45    |
|                | ACT     | 339   | 10     | 55           | 35    | 69      | 61    |
|                | CD      | 178   | 10     | 97           | 74    | 94      | 89    |
| Light-GBM      | GC      | 92    | 10     | 68           | 84    | 31      | 59    |
|                | AF      | 325   | 10     | 30           | 55    | 42      | 24    |
|                | BH      | 118   | 10     | 15           | 17    | 16      | 25    |
|                | BS      | 287   | 10     | 119          | 48    | 68      | 116   |
|                | WQR     | 454   | 10     | 64           | 29    | 129     | 128   |
|                | ACT     | 132   | 10     | 54           | 46    | 16      | 51    |
| Mean           | CD      | 409   | 10     | 68           | 53    | 12      | 50    |
|                | GC      | 501   | 10     | 61           | 86    | 16      | 35    |
| Mean           |         | 223   | 10     | 64           | 60    | 60      | 66    |

1311  
1312  
1313  
1314  
1315  
1316  
1317  
1318  
1319  
1320  
1321  
1322  
1323  
1324  
1325  
1326  
1327  
1328  
1329  
1330  
1331  
1332  
1333  
1334  
1335  
1336  
1337  
1338  
1339  
1340  
1341  
1342  
1343  
1344  
1345  
1346  
1347  
1348  
1349

## D.10 COMPUTATION COST

Table 23 compares the number of features evaluated during the feature search process. Guided by domain knowledge, our method FEBP evaluates much fewer features than DIFER (Zhu et al., 2022b) and OpenFE (Zhang et al., 2023).

Tables 24 and 25 summarize the computation time, with *gpt-3.5-turbo-0125* as the LLM. For FEBP, the computation time of LLM generation and feature evaluation is relatively stable across datasets of varying sizes. We note that the LLM generation time can be substantially reduced by instructing the LLM to generate multiple features in a generation step.

**Table 23:** Comparison of the number of evaluated features during feature search.

| Model          | Dataset | DIFER | OpenFE | FEBP |
|----------------|---------|-------|--------|------|
| Linear Model   | AF      | 2083  | 224    | 200  |
|                | BH      | 2081  | 1167   | 200  |
|                | WQR     | 2083  | 929    | 200  |
|                | ACT     | 2077  | 4310   | 200  |
|                | CD      | 2088  | 3385   | 200  |
|                | GC      | 2076  | 4169   | 200  |
| Random Forests | AF      | 2085  | 224    | 200  |
|                | BH      | 2079  | 1051   | 200  |
|                | BS      | 2082  | 310    | 200  |
|                | WQR     | 2085  | 929    | 200  |
|                | ACT     | 2079  | 1636   | 200  |
|                | CD      | 2086  | 1801   | 200  |
| Light-GBM      | GC      | 2078  | 2139   | 200  |
|                | AF      | 2084  | 224    | 200  |
|                | BH      | 2080  | 1051   | 200  |
|                | BS      | 2083  | 310    | 200  |
|                | WQR     | 2084  | 929    | 200  |
|                | ACT     | 2079  | 1636   | 200  |
| Mean           | CD      | 2087  | 1801   | 200  |
|                | GC      | 2078  | 2139   | 200  |

## D.11 POTENTIAL FAILURE MODES

One potential failure mode is the generation of features that are duplicates of existing candidate features or syntactically invalid. The third column of each configuration in Table 3 reports the number of LLM responses needed to construct 200 candidate features in our experiments. Specifically, the proportion of valid new features is around 55% using GPT-3.5 and around 60% using GPT-4 on average. Feature search tends to converge per our feature divergence analysis in Section 5.6 and could get stuck in local optima when example features in the prompt are highly similar. From Figure 10, the number of LLM responses needed to construct a new candidate feature remains relatively stable as the algorithm iterates, suggesting low likelihood of getting stuck. Table 4 shows that including more example features in the prompt improves the success rate of feature construction on average by increasing the diversity. Another potential failure mode is that the generated explanation of a constructed feature may be inaccurate, e.g., the column index may be inconsistent with the feature name.

## E EXPERIMENTS ON PROPRIETARY DATASETS

We have conducted experiments on our proprietary real-world dataset containing over 100,000 samples and over 1,000 features, with most features contain a substantial proportion of missing values. We perform AutoFE on features of top 100 mutual information scores with the target. PromptFE brings significant performance improvements to downstream models on our proprietary dataset.

1404  
 1405  
 1406  
 1407  
 1408  
 1409  
 1410  
 1411  
 1412  
 1413  
 1414  
 1415  
 1416  
 1417  
 1418  
 1419  
 1420  
 1421  
 1422  
 1423  
 1424  
 1425  
 1426  
 1427  
 1428  
 1429  
 1430  
 1431  
 1432  
 1433  
 1434  
 1435  
 1436  
 1437  
 1438  
 1439  
 1440  
 1441  
 1442  
 1443  
 1444  
 1445  
 1446  
 1447  
 1448  
 1449  
 1450  
 1451  
 1452  
 1453  
 1454  
 1455  
 1456  
 1457

**Table 24:** Comparison of computation time, in minutes.

| Model          | Dataset | DIFER  | OpenFE | CAAFE | FEBP  |
|----------------|---------|--------|--------|-------|-------|
| Linear Model   | AF      | 33.49  | 0.21   | 1.73  | 42.80 |
|                | BH      | 41.17  | 0.21   | 1.18  | 41.28 |
|                | WQR     | 34.94  | 0.25   | 1.21  | 42.33 |
|                | ACT     | 44.18  | 0.40   | 1.25  | 43.60 |
|                | CD      | 433.94 | 1.49   | 3.17  | 57.82 |
|                | GC      | 29.30  | 0.37   | 1.68  | 43.71 |
| Random Forests | AF      | 178.50 | 0.23   | 4.22  | 63.30 |
|                | BH      | 89.07  | 0.24   | 5.52  | 51.70 |
|                | BS      | 98.50  | 0.23   | 4.05  | 51.13 |
|                | WQR     | 298.46 | 0.29   | 9.35  | 63.12 |
|                | ACT     | 78.44  | 0.28   | 3.82  | 44.66 |
|                | CD      | 571.33 | 1.12   | 14.05 | 94.08 |
| Light-GBM      | GC      | 60.41  | 0.28   | 3.24  | 45.06 |
|                | AF      | 301.56 | 0.25   | 5.81  | 63.06 |
|                | BH      | 62.30  | 0.24   | 3.01  | 44.84 |
|                | BS      | 74.59  | 0.24   | 2.55  | 45.23 |
|                | WQR     | 361.19 | 0.29   | 5.68  | 58.97 |
|                | ACT     | 36.39  | 0.28   | 1.73  | 42.71 |
| Mean           | CD      | 102.04 | 1.07   | 2.49  | 46.34 |
|                | GC      | 48.63  | 0.28   | 2.97  | 43.03 |
| Mean           |         | 148.92 | 0.41   | 3.94  | 51.44 |

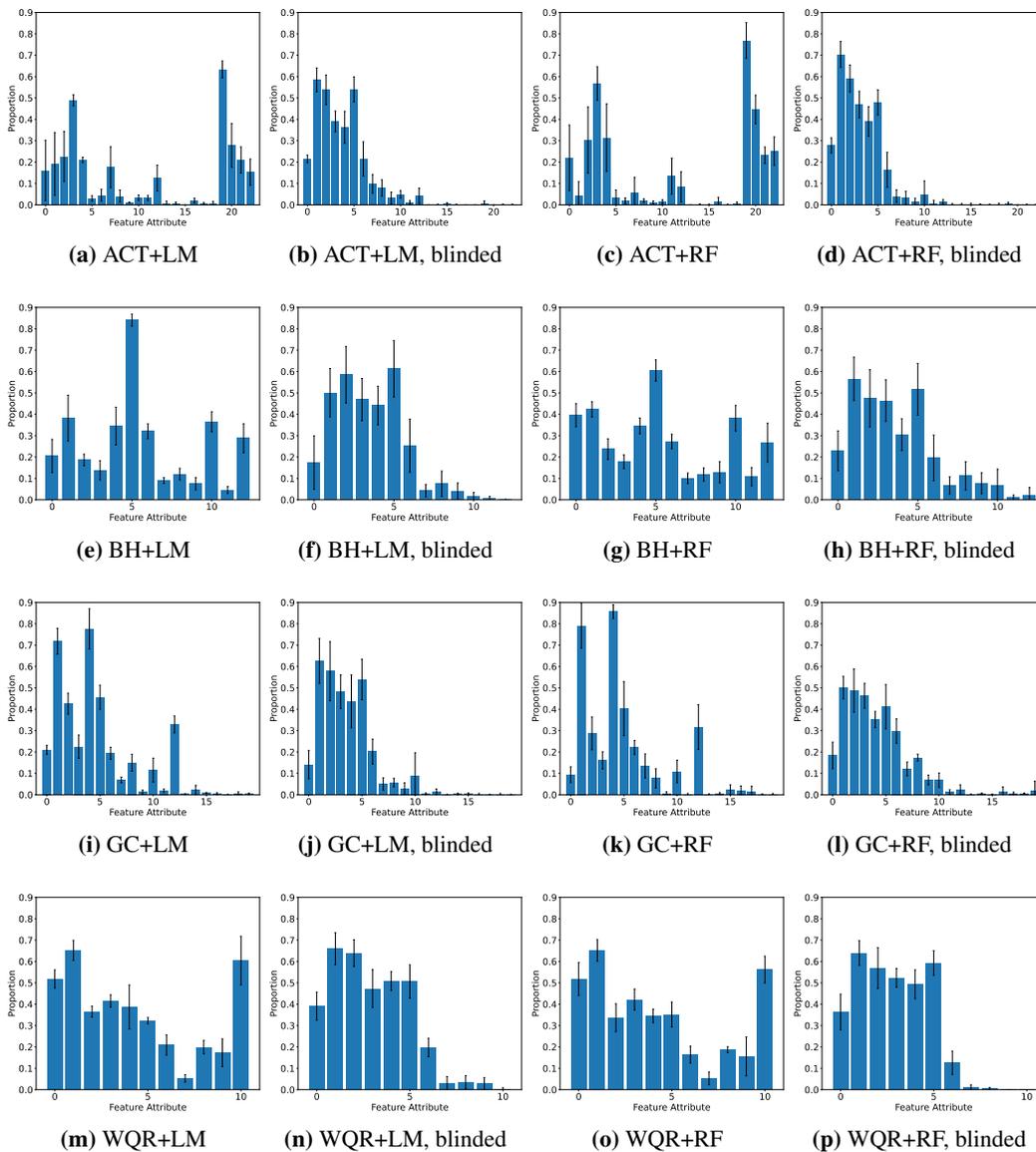
**Table 25:** Computation time of different components of FEBP, in minutes.

| Model          | Dataset | LLM Generation | Feature Evaluation | Feature Selection |
|----------------|---------|----------------|--------------------|-------------------|
| Linear Model   | AF      | 16.73          | 22.98              | 3.08              |
|                | BH      | 18.50          | 20.18              | 2.60              |
|                | WQR     | 19.07          | 20.24              | 3.02              |
|                | ACT     | 18.92          | 20.97              | 3.71              |
|                | CD      | 16.73          | 25.14              | 15.95             |
|                | GC      | 17.01          | 23.24              | 3.47              |
| Random Forests | AF      | 15.34          | 25.32              | 22.64             |
|                | BH      | 18.60          | 23.69              | 9.41              |
|                | BS      | 15.12          | 25.16              | 10.87             |
|                | WQR     | 12.75          | 23.81              | 26.56             |
|                | ACT     | 13.79          | 21.67              | 9.20              |
|                | CD      | 12.48          | 25.89              | 55.71             |
| Light-GBM      | GC      | 14.80          | 21.91              | 8.35              |
|                | AF      | 17.37          | 21.06              | 24.63             |
|                | BH      | 19.70          | 20.40              | 4.74              |
|                | BS      | 17.03          | 22.18              | 6.02              |
|                | WQR     | 16.27          | 21.19              | 21.51             |
|                | ACT     | 19.18          | 20.24              | 3.29              |
| Mean           | CD      | 16.53          | 21.68              | 8.13              |
|                | GC      | 17.00          | 20.40              | 5.63              |
| Mean           |         | 16.65          | 22.37              | 12.43             |

## F ADDITIONAL ANALYSIS

### F.1 FEATURE ANALYSIS

Figure 15 compares the proportions of generated features selecting each feature attribute across different datasets and downstream models (linear models and Random Forests) for both the full and semantically blinded versions of FEBP. In the blinded version, we observe that the LLM tends to prioritize earlier feature attributes in the dataset while paying less attention to later ones, reflecting an inherent bias of the language model. In contrast, in the full version, the selection of feature attributes is guided by the semantic information of the dataset rather than the positional order of the attributes. Specifically, Attribute 19 *CD4 at baseline* in AIDS Clinical Trials (ACT) and Attribute 10 *alcohol* in Wine Quality Red (WQR), which contain useful information for predicting the targets *censoring indicator* and *quality*, respectively, are included in the majority of the generated features. This demonstrates the role of dataset semantic information in the LLM-based feature search process.

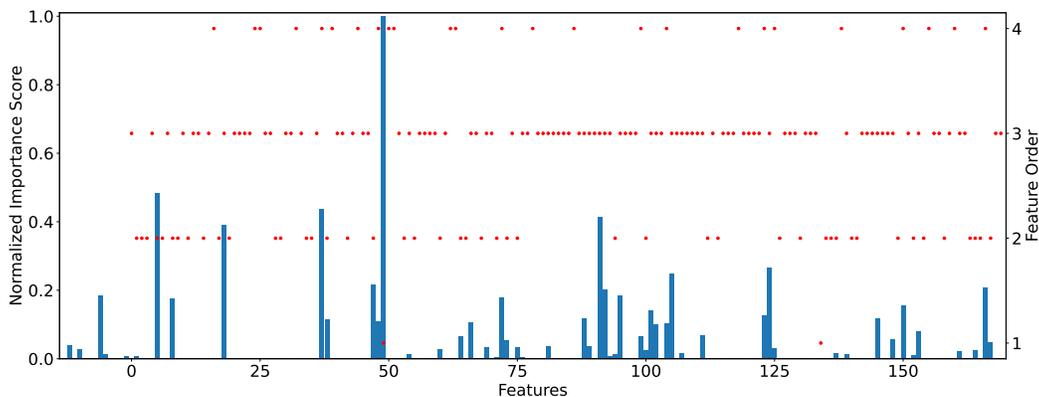


**Figure 15:** The proportions of generated features selecting each feature attribute in the dataset.

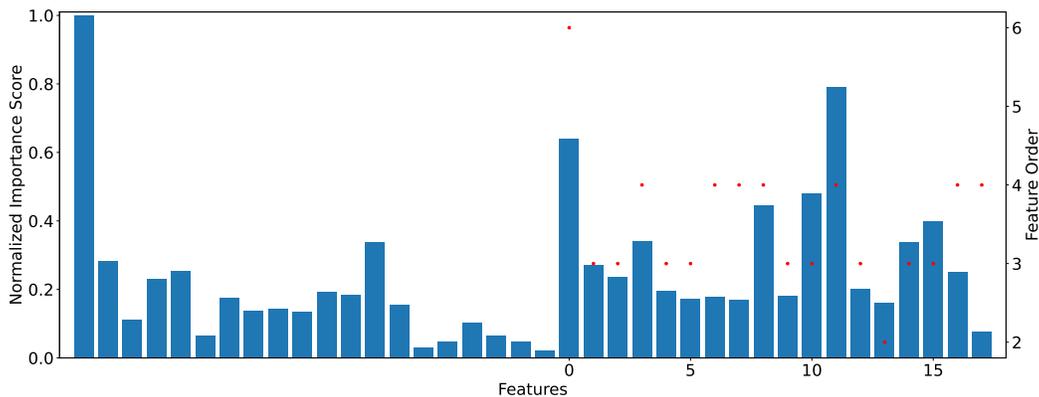
1512  
1513  
1514  
1515  
1516  
1517  
1518  
1519  
1520  
1521  
1522  
1523  
1524  
1525  
1526  
1527  
1528  
1529  
1530  
1531  
1532  
1533  
1534  
1535  
1536  
1537  
1538  
1539  
1540  
1541  
1542  
1543  
1544  
1545  
1546  
1547  
1548  
1549  
1550  
1551  
1552  
1553  
1554  
1555  
1556  
1557  
1558  
1559  
1560  
1561  
1562  
1563  
1564  
1565

## F.2 FEATURE IMPORTANCE

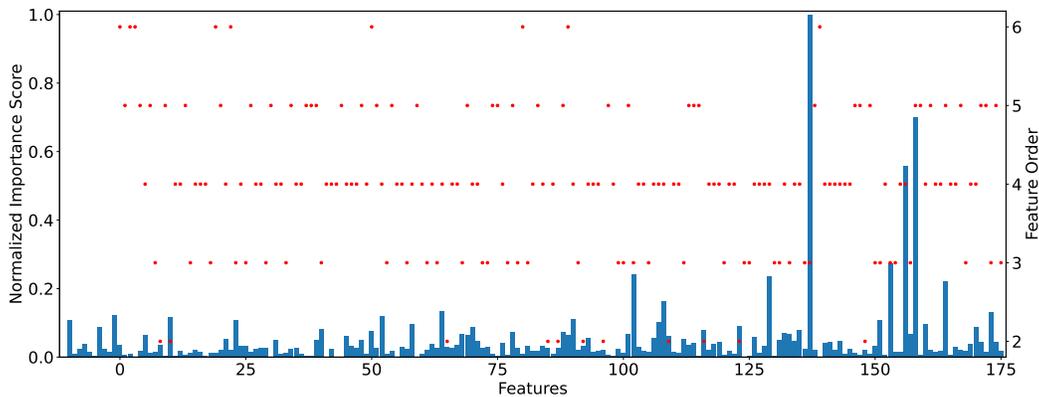
Figure 16 shows the feature importance across different datasets and downstream models. For linear models, we use the magnitudes of coefficients; for Random Forests (Breiman, 2001), we use the impurity-based feature importance; for LightGBM (Ke et al., 2017), we use the total gains of splits. FEBP enhances the datasets with generated features that extract valuable information for target prediction. Additionally, we observe that Random Forests and LightGBM benefit from features of higher orders compared to linear models, as they are capable of synthesizing simple features internally. Our method adaptively adjusts the feature complexity to suit different downstream models.



(a) BH+LM. Test performance improves from to 0.3776 to 0.5157.



(b) GC+RF. Test performance improves from to 0.7450 to 0.7700.



(c) WQR+LGBM. Test performance improves from to 0.3825 to 0.4299.

**Figure 16:** The blue bars show the normalized feature importance scores. The red dots show the order of features. The generated features are positioned on the x-axis starting at index 0, following the raw features.

1566 F.3 FEATURE SEARCH EXAMPLE  
1567

1568 Figure 17 shows samples of the LLM-based feature search process from a single run. In iteration  
1569 155, the generated feature is derived from the combination of three example features. In iteration  
1570 158, the generated feature is derived from the combination of two example features, with the fea-  
1571 ture attributes modified. Iteration 160 exhibits an instance of modification on the example feature.  
1572 Iteration 161 exhibits an instance of analogy of the example feature.  
1573

1574 Figure 17: Examples of feature search. The prompt examples that relate to the generated feature  
1575 are highlighted in red.

1576  
1577 Iteration 155  
1578 Generated feature:  
1579 6, 12, division, 18, \*, min\_max  
1580 Prompt examples:  
1581 4, log, 13, \*, 17, min\_max, sqrt\_abs, +  
1582 6, 12, \*, 21, sqrt\_abs, +  
1583 4, 14, +, 18, \*, min\_max  
1584 11, 21, +, 6, \*, sqrt\_abs  
1585 2, 4, +, 6, division, log  
1586 4, 6, +, 16, \*, 17, +, min\_max  
1587 6, 21, division, 17, log, \*  
1588 12, min\_max, 18, \*, 5, mod\_column  
1589 2, 6, mod\_column, 17, +, sqrt\_abs  
1590 2, 3, +, 4, sqrt\_abs, \*  
1591 4, 12, division, 16, \*, min\_max, sqrt\_abs  
1592  
1593 Iteration 158  
1594 Generated feature:  
1595 5, log, 19, \*, 15, +, min\_max  
1596 Prompt examples:  
1597 4, log, 13, \*, 17, min\_max, sqrt\_abs, +  
1598 6, 12, \*, 21, sqrt\_abs, +  
1599 4, 14, +, 18, \*, min\_max  
1600 11, 21, +, 6, \*, sqrt\_abs  
1601 2, 4, +, 6, division, log  
1602 4, 6, +, 16, \*, 17, +, min\_max  
1603 6, 21, division, 17, log, \*  
1604 12, min\_max, 18, \*, 5, mod\_column  
1605 2, 6, mod\_column, 17, +, sqrt\_abs  
1606 2, 3, +, 4, sqrt\_abs, \*  
1607 4, 12, division, 16, \*, min\_max, sqrt\_abs  
1608  
1609 Iteration 160  
1610 Generated feature:  
1611 11, 21, \*, 4, +, 14, mod\_column  
1612 Prompt examples:  
1613 4, log, 13, \*, 17, min\_max, sqrt\_abs, +  
1614 6, 12, \*, 21, sqrt\_abs, +  
1615 4, 14, +, 18, \*, min\_max  
1616 11, 21, +, 6, \*, sqrt\_abs  
1617 2, 4, +, 6, division, log  
1618 4, 6, +, 16, \*, 17, +, min\_max  
1619 6, 21, division, 17, log, \*  
12, min\_max, 18, \*, 5, mod\_column  
2, 6, mod\_column, 17, +, sqrt\_abs  
2, 3, +, 4, sqrt\_abs, \*  
4, 12, division, 16, \*, min\_max, sqrt\_abs  
Iteration 161  
Generated feature:  
6, 12, +, 17, \*, min\_max  
Prompt examples:

1620  
1621  
1622  
1623  
1624  
1625  
1626  
1627  
1628  
1629  
1630  
1631  
1632  
1633  
1634  
1635  
1636  
1637  
1638  
1639  
1640  
1641  
1642  
1643  
1644  
1645  
1646  
1647  
1648  
1649  
1650  
1651  
1652  
1653  
1654  
1655  
1656  
1657  
1658  
1659  
1660  
1661  
1662  
1663  
1664  
1665  
1666  
1667  
1668  
1669  
1670  
1671  
1672  
1673

```

6, 12, *, 21, sqrt_abs, +
4, 14, +, 18, *, min_max
11, 21, *, 4, +, 14, mod_column
11, 21, +, 6, *, sqrt_abs
2, 4, +, 6, division, log
4, 6, +, 16, *, 17, +, min_max
6, 21, division, 17, log, *
12, min_max, 18, *, 5, mod_column
2, 6, mod_column, 17, +, sqrt_abs
2, 3, +, 4, sqrt_abs, *
4, 12, division, 16, *, min_max, sqrt_abs

```

## G FURTHER DISCUSSION ON DIFFERENCES FROM EXISTING WORKS

Although our work FEBP and CAAFE (Hollmann et al., 2023) both utilize LLMs to construct new features incorporating dataset semantic information, they differ in several key aspects. We design FEBP such that it taps into the in-context learning capability of LLMs and performs effective feature search. In FEBP, we provide top-performing constructed features in the prompt as learning examples, label them with performance scores, and rank them by score. We demonstrate that the LLM learns to optimize feature construction over the course of algorithm. CAAFE instead stores all previous instructions and code snippets in the conversation history, which hinders the in-context learning of optimal feature patterns. It quickly consumes the LLM’s context as the algorithm iterates, incurring more and more LLM generation costs. In comparison, the LLM generation cost of FEBP stays constant across iterations, without a maximum limit on the number of iterations it can perform. Therefore, our method FEBP has stronger capability of performing feature search in large search spaces requiring many iterations, such as datasets with numerous feature attributes.

In FEBP, we also explore representing features in a different form, i.e., canonical RPN (cRPN). We refer to Appendix A for further detail. Compared with the Python code representation in CAAFE, cRPN is more compact, which not only reduces LLM generation costs but also makes the in-context learning of feature patterns easier, and more human interpretable. The use of pre-defined operators reduces the search space and simplifies the learning process for optimizing feature construction. Together, our approach gives better control than code representation that helps avoid undesirable or unexpected LLM outputs. Another advantage of cRPN is that it is convenient to import external features (as outlined in Algorithm 1) and export the results as individual features, providing compatibility with other feature engineering methods.

More fundamentally, we demonstrate in this work that general-purpose LLMs like GPTs can effectively model recursive tree structures in the form of cRPN feature expressions and reason about the structures in the context of semantic information, shedding light on further LLM-driven applications. We hereby underscore the importance of adopting proper representation for the downstream task to tap into LLMs’ potential.

FeatLLM, LFG