000 ITERATIVE FEATURE SPACE OPTIMIZATION THROUGH 001 INCREMENTAL ADAPTIVE EVALUATION 002 003

Anonymous authors

Paper under double-blind review

ABSTRACT

Iterative feature space optimization involves systematically evaluating and adjust-012 ing the feature space to improve downstream task performance. However, existing works suffer from three key limitations: 1) overlooking differences among data samples leads to evaluation bias; 2) tailoring feature spaces to specific machine learning models results in overfitting and poor generalization; 3) requiring 015 the evaluator to be retrained from scratch during each optimization iteration sig-016 nificantly reduces the overall efficiency of the optimization process. To bridge these gaps, we propose a <u>gEneralized</u> <u>A</u>daptive feature <u>Space</u> <u>E</u>valuator (EASE) 018 to efficiently produce optimal and generalized feature spaces. This framework 019 consists of two key components: Feature-Sample Subspace Generator and Contextual Attention Evaluator. The first component aims to decouple the information distribution within the feature space to mitigate evaluation bias. To achieve this, we first identify features most relevant to prediction tasks and samples most challenging for evaluation based on feedback from the subsequent evaluator. These 024 identified feature and samples are then used to construct feature subspaces for next optimization iteration. This decoupling strategy makes the evaluator consistently target the most challenging aspects of the feature space. The second component intends to incrementally capture evolving patterns of the feature space for efficient evaluation. We propose a weighted-sharing multi-head attention mechanism 028 to encode key characteristics of the feature space into an embedding vector for evaluation. Moreover, the evaluator is updated incrementally, retaining prior evaluation knowledge while incorporating new insights, as consecutive feature spaces during the optimization process share partial information. Extensive experiments on twelve real-world datasets demonstrate the effectiveness of the proposed framework. Our code and data are publicly available ¹.

034

037

033

004

006

008 009

010 011

013

014

017

021

025

026

027

029

031

INTRODUCTION 1

038 Iterative feature space optimization systematically evaluates and refines the feature space to enhance downstream task performance (Jia et al., 2022). As depicted in Figure 1a, the optimization module 040 iteratively enhances the feature space based on the feedback from the evaluator. This optimization 041 process continues until the optimal feature space is identified. This approach has demonstrated 042 broad applicability and has been successfully adopted in various fields, including biology, finance, 043 and medicine (Zhu et al., 2023; Htun et al., 2023; Vommi & Battula, 2023).

044 Research in this domain has received significant attention (Zebari et al., 2020). Recursive optimiza-045 tion methods focus on evaluating feature importance to progressively refine the feature space (Darst 046 et al., 2018; Priyatno et al., 2024). For instance, Escanilla et al. (2018) utilized sensitivity testing 047 with membership queries on trained models to recursively identify key features. To improve the ef-048 ficiency of these recursive methods, evolutionary algorithms and reinforcement learning (RL) were subsequently introduced, further accelerating the refinement process (Xiao et al., 2024; Wang et al., 2024; Liu et al., 2021a). For example, Wang et al. (2022) employed three cascading agents to replicate the feature engineering process typically performed by human experts, using RL to streamline 051 the exploration phase. 052

¹Https://anonymous.4open.science/r/EASE-1C51

081

082

084

085

087



Figure 1: (a) Illustration of the iterative feature space optimization, where the optimization module
 refines the feature space based on the feedback of the evaluator until the optimal one is identified.
 (b) The feature spaces between consecutive iterations exhibit informational overlap.

066 But, existing approaches suffer from three key limitations: Limitation 1: Evaluation bias. These 067 methods do not account for variability between samples, which limits the evaluator's ability to cap-068 ture the full range of features of the space. As a result, the performance assessments become biased 069 and do not objectively reflect the quality of the feature space. Limitation 2: Non-generalizability. Customizing the feature space by replacing the evaluator based on specific requirements limits its 071 ability to capture generalizable patterns. Consequently, the resulting feature space lacks flexibility 072 and cannot be effectively applied across diverse scenarios. Limitation 3: Training Inefficiency. Retraining the evaluator from scratch at each iteration significantly increases computational de-073 mands. This technical trait leads to a time-consuming process that hinders efficiency and scalability. 074

Thus, there is a vital need for a robust evaluation framework that can efficiently assess feature space quality, enabling the creation of generalized and optimal feature spaces. This framework should integrate seamlessly with iterative feature space optimization algorithms to enhance their performance and efficiency. However, to accomplish this, there are two key technical challenges:

- Challenge 1: Complicated Feature Interactions. Within the feature space, complex feature-feature interactions exist, which are important for understanding its characteristics and enabling more effective refinement. But, how can we effectively capture such complicated information as the guidance during the iterative optimization process?
- Challenge 2: Incremental Evaluator Updates. As illustrated in Figure 1b, the feature spaces between consecutive iterations exhibit partial overlap. This overlap presents an opportunity to update the evaluator efficiently, rather than retraining it from scratch. But, how can we incrementally update the parameters of the evaluator to ensure it retains essential prior evaluation knowledge while simultaneously integrating new evaluation insights?

To address these challenges, we propose **EASE**, a <u>gEneralized</u> <u>A</u>daptive feature <u>Space</u> <u>E</u>valuator, 090 which can seamlessly integrate as a plugin into any iterative feature space optimization method for 091 enhanced feature space refinement. This framework contains two key components: Feature-Sample 092 Subspace Generator and Contextual Attention Evaluator. The first component aims to decouple the information distribution within the feature space to mitigate evaluation bias. To achieve this, we initially employ the feature index optimizer to select the features most relevant to the prediction task. 094 Next, we use the sample index optimizer to identify the samples that present the greatest evaluation 095 challenges. Both the previous two steps were guided by feedback from the subsequent evaluator. Fi-096 nally, we use the identified features and samples to construct feature subspaces for the next iteration of feature space refinement. This decoupling strategy enables the evaluator to consistently target the 098 most challenging aspects of the feature space, thereby facilitating the comprehensive comprehension and establishing a robust foundation for objective evaluation. The second component intends 100 to incrementally capture the evolving patterns within the feature space for enhancing evaluation ef-101 ficiency. Specifically, we employ a multi-head attention mechanism as the backbone to develop the 102 evaluator. The feature subspaces are sequentially fed into the evaluator to capture complex rela-103 tionships, leveraging contextual information across subspaces. The evaluator utilizes shared model 104 weights across various feature subspaces. Moreover, since refined feature spaces across consecutive 105 optimization iterations often share overlapping information, we incrementally update the evaluator's parameters to retain prior evaluation knowledge while incorporating new insights from the evolv-106 ing feature space. Finally, we apply EASE to iterative feature selection algorithms and conduct 107 extensive experiments on twelve real-world datasets to validate its superiority and effectiveness.

108 2 RELATED WORK

109 110

Incremental Learning (IL). IL aims to acquire new knowledge without forgetting the knowledge 111 it has already learned (Zhu et al., 2021). IL is applied in scenarios such as dynamic environments 112 (Shieh et al., 2020; Read et al., 2012) and online learning (Shim et al., 2021). IL methods can 113 be divided into three categories: regularization, memory replay, and parameter isolation methods. 114 Regularization-based methods (e.g., Kirkpatrick et al. (2017); Li & Hoiem (2017)) prevent significant changes in important parameters of previous tasks. Memory replay methods retain old task data 115 116 (Isele & Cosgun, 2018) or use generative models to simulate it (Shin et al., 2017), and then train this data alongside new task data when learning new tasks. Parameter isolation methods achieve 117 task isolation by assigning independent model parameters to different tasks (e.g., Rajasegaran et al. 118 (2019); Serra et al. (2018)) or by expanding the network structure to accommodate new tasks (e.g., 119 Moriya et al. (2018); Aljundi et al. (2017)). In this paper, we update the evaluator using the Elastic 120 Weight Consolidation (EWC) strategy (Kirkpatrick et al., 2017; Liu et al., 2021b). This approach 121 estimates the importance of model parameters for previous tasks and minimizes changes to these 122 important parameters when training on new tasks. This method significantly improves the training 123 efficiency. 124

Multi-Head Attention. The multi-head attention mechanism enhances representation capability by simultaneously attending to different subspaces of the input data (Vaswani, 2017; Messaoud et al., 2021). This technique is used in natural language processing (Vaswani, 2017; Sun et al., 2020) and object detection (Dai et al., 2021) to capture complex patterns and dependencies. Unlike previous works, we propose a weighted multi-head attention mechanism that shares weights to encode key characteristics of the feature space into an embedding vector for the evaluation.

Feature Selection (FS). FS is widely used in high-dimensional fields (Nguyen et al., 2020), such 131 as bioinformatics (Pudjihartono et al., 2022) and finance (Arora & Kaur, 2020). Among these FS 132 methods, the wrapper method stands out for its ability to select features based directly on model 133 performance (Nouri-Moghaddam et al., 2021). Wrapper-based methods use the performance of the 134 downstream model as a criterion and employ iterative search to find the optimal feature subset (Liu 135 et al., 2023). The most representative wrapper method is Recursive Feature Elimination (RFE). 136 RFE iteratively trains the model and removes the least important features, gradually reducing the 137 feature set size until a specified criterion is met (Guyon et al., 2002). In this paper, we use FS as a 138 representative example of feature optimization to illustrate the subsequent technical details.

139 140 141

142

149

150

155 156 157

3 PROBLEM STATEMENT

This paper introduces a novel feature space evaluator to efficiently identify the optimal feature space. The proposed evaluator can be seamlessly integrated into any iterative feature space optimization algorithm. Formally, given a dataset $\mathbb{D} = \langle \mathcal{F}, y \rangle$, where \mathcal{F} represents the feature space and ydenotes the target label space, we first initialize the parameters $\Theta_{\mathcal{M}}$ of the evaluator \mathcal{M} based on \mathbb{D} . This initialization is achieved by minimizing the prediction error \mathcal{L} . The learning objective can be defined as:

$$\arg\min_{\boldsymbol{\Theta}_{\mathcal{M}}} \mathcal{L}(\mathcal{M}(\boldsymbol{\mathcal{F}};\boldsymbol{\Theta}_{\mathcal{M}}),\boldsymbol{y}).$$
(1)

In the *t*-th optimization, we can get a new feature space $\langle \mathcal{F}^{(t)}, y^{(t)} \rangle$. We leverage the information overlap between feature spaces from consecutive iterations to incrementally update $\Theta_{\mathcal{M}}^{(t)}$, enabling efficient tracking of evolving patterns and providing an accurate evaluation of the feature space. The learning process can be formulated as follows:

$$\arg\min_{\boldsymbol{\Theta}_{\mathcal{M}}} \mathcal{L}(\mathcal{M}(\boldsymbol{\mathcal{F}}^{(t)};\boldsymbol{\Theta}_{\mathcal{M}}^{(t)}),\boldsymbol{y}^{(t)}) + \lambda \|\boldsymbol{\Theta}_{\mathcal{M}}^{(t)} - \boldsymbol{\Theta}_{\mathcal{M}}^{(t-1)}\|_{2},$$
(2)

where λ is a regularization parameter that balances retaining prior evaluation knowledge with incorporating new insights from the updated feature space, and $\|\cdot\|_2$ is L2 norm. The learning process continues until either the maximum number of iterations is reached or the optimal feature space is identified. The design and optimization of \mathcal{M} represent the core contribution of this paper. For clarity, key notations are summarized in Table 4 in the Appendix.



Figure 2: Framework overview and parameter update for EASE. The framework comprises two key components: the Feature-Sample Subspace Generator and the Contextual Attention Evaluator. The first component aims to decouple the complex information within the feature space, enabling the evaluator to focus on capturing the most challenging aspects for evaluation. The second component is designed to comprehensively capture the characteristics of the feature space, ensuring fair and accurate evaluation. By considering information overlap across consecutive iterations, the evaluator incrementally updates its parameters, enhancing the efficiency of the overall optimization process.

185

4 Methodology

186 Framework Overview. Figure 2a shows the framework overview of EASE, which includes two key 187 components: 1) Feature-Sample Subspace Generator; and 2) Contextual Attention Evaluator. The 188 first component decouples the information distribution within the evolving feature space, aiming to 189 reduce the complexity for the subsequent evaluator in capturing the key characteristics of the feature 190 space. Specifically, in each optimization, guided by the evaluator, we first use the feature index op-191 timizer to identify the most important features for the downstream prediction task. Then, the sample index optimizer is used to discover samples that are challenging to evaluate. After that, we construct 192 fixed-length feature subspaces by applying a random combination strategy to the identified features 193 and samples. The second component aims to efficiently capture complex feature interactions within 194 the feature space to facilitate effective evaluation. In detail, we first input the multiple feature sub-195 spaces constructed by the first component into the contextual attention evaluator to capture complex 196 interactions and encode them into an embedding vector. In addition, the embedding vector is used to 197 perform the evaluation task and the resulting prediction error is fed back to the previous component as guidance. During this process, as illustrated in Figure 2b, considering the partial information 199 overlap between feature spaces in consecutive iterations, we incrementally update the evaluator's 200 parameters to efficiently incorporate new evaluation insights.

201 202

203 204

4.1 FEATURE-SAMPLE SUBSPACE GENERATOR

Why is feature space decoupling important? During the iterative feature space optimization process, complex feature interactions can obscure the underlying patterns. A comprehensive understanding of these interactions is crucial for an accurate and objective evaluation. By decoupling the feature space, we can reduce the complexity of the learning task, allowing the evaluator to concentrate on the most challenging aspects, ultimately resulting in more effective and precise evaluations.

Feature Index Optimizer is designed to identify the features most relevant to the subsequent evaluation task. We derive the feature index subset based on the feature importance scores. Formally, in the *t*-th iteration, given the feature space $\mathcal{F}^{(t)}$ and the target variable $y^{(t)}$, the importance score Score (f_i) for each feature f_i is calculated by assessing the impact of removing that feature on the performance of the model. The importance score is computed as follows:

215

$$Score(\boldsymbol{f}_i) = \mathcal{L}(\mathcal{M}(\boldsymbol{\mathcal{F}}^{(t)}; \boldsymbol{\Theta}_{\mathcal{M}}^{(t)})) - \mathcal{L}(\mathcal{M}(\boldsymbol{\mathcal{F}}^{(t)} \setminus \{\boldsymbol{f}_i\}; \boldsymbol{\Theta}_{\mathcal{M}}^{(t)})).$$
(3)

216 Here, $\mathcal{L}(\mathcal{M}(\cdot; \Theta_{\mathcal{M}}^{(t)})))$ denotes the feature space evaluator that measures model performance, and $\mathcal{F}^{(t)} \setminus \{f_i\}$ represents the feature space $\mathcal{F}^{(t)}$ with feature f_i omitted. This means that we can esti-217 218 mate the Score(f_i) by removing f_i and observing the change in the loss. The loss function can be 219 tailored to the specific task. For classification, cross-entropy loss is commonly used, whereas for re-220 gression, mean squared error is typically employed. Once the importance scores for all features have 221 been computed, they are ranked in descending order. The features with higher ranks are considered 222 to be the most significant contributors to the performance of the model. We select the best k features based on their importance scores to identify the subset of the feature index $f^{(t)} = \{f_1, f_2, \cdots, f_k\}$ 224 that is the most relevant to the evaluation task. This component can be replaced with any feature selection module, allowing EASE to be compatible with any iterative feature selection algorithm. 225

226 Sample Index Optimizer aims to select the most challenging samples for the subsequent evalua-227 tion task. The sample index subset is derived based on the evaluation error from the feature space 228 evaluator. Formally, in the t-th iteration, the feature space consists of n samples and a target variable **y**. For the *i*-th sample x_i , the prediction error is given by $\mathcal{L}_i^{(t-1)} = \ell(\mathbf{y}_i, \hat{\mathbf{y}}_i)$, where ℓ denotes the evaluation metric, \mathbf{y}_i is the target value, and $\hat{\mathbf{y}}_i$ is the prediction value. The sampling probability for sample x_i is defined as: $P(X = x_i) = \frac{\mathcal{L}_i^{(t-1)}}{\sum_{j=1}^n \mathcal{L}_j^{(t-1)}}$, where $P(X = x_i)$ represents the likelihood of calcoring complex is based on its relative representation. 229 230 231 232 of selecting sample x_i based on its relative prediction error. This ensures that samples with large 233 errors have a higher probability of being selected. To efficiently sample from this distribution, we 234 use the cumulative distribution function (CDF), which allows us to transform a uniformly distributed 235 random number into a sample from the desired probability distribution. The CDF is constructed as: 236 $C_i = \sum_{k=1}^{i} p_k = \sum_{k=1}^{i} P(X = x_k)$, where C_i represents the cumulative sum of probabilities up 237 to the i-th sample. More specifically, in the weighted sampling process, we first generate a random 238 number r uniformly distributed in the interval [0, 1]. Next, we identify the first sample index i such 239 that the CDF satisfies $C_i \ge r$, and select the corresponding sample x_i . The process is repeated until 240 a new set of sample indices $\mathbb{I}^{(t)}$ is collected. Using the CDF, the sampling process aligns with the 241 distribution of prediction errors, giving higher priority to samples with larger errors. 242

Feature Subspace Construction decouples the complex feature space into distinct portions to im-243 prove the understanding of the subsequent evaluator for a fair evaluation. Thus, we introduce the 244 strategy for the construction of feature subspaces using the set of feature index $f^{(t)}$ and sample 245 index $\mathbb{I}^{(t)}$. More specifically, we sample s sample indices from $\mathbb{I}^{(t)}$ for q times using repeated sam-246 pling to obtain various sub sample indices $\{\mathbb{I}_1^{(t)}, \mathbb{I}_2^{(t)}, \cdots, \mathbb{I}_q^{(t)}\}$ with the same length. Next, we use 247 the obtained sample indices and $f^{(t)}$ to select the corresponding samples and features, constructing 248 various feature subspaces denoted as $\mathcal{B}^{(t)} = \{B_1^{(t)}, B_2^{(t)}, \cdots, B_q^{(t)}\}$. The *i*-th feature subspace 249 250 $B_i^{(t)} \in \mathbb{R}^{s \times k}$, where s is the number of data samples and k is the number of features within $B_i^{(t)}$. 251 Through this process, we decouple the information distribution within the feature space, preserv-252 ing the most important and challenging aspects for evaluation in $\mathcal{B}^{(t)}$. The pseudo-code for feature 253 subspace construction is provided in Algorithm 1 in the Appendix for improving reproducibility. 254

- 255 4.2 CONTEXTUAL
- 256

265

4.2 CONTEXTUAL ATTENTION EVALUATOR

To thoroughly capture the complicated interactions of the feature space, we design a contextual attention evaluator leveraging a multi-attention mechanism. The learned feature subspaces $\mathcal{B}^{(t)}$, are sequentially fed into the evaluator to facilitate comprehensive information extraction. We use the *i*-th feature subspace to illustrate the following calculation process. For clarity, we omit the (t)notation for the *i*-th example.

More specifically, we begin by projecting B_i into three different spaces: the query Q_i , key K_i , and value V_i spaces. These projections are computed through linear transformations, which can be defined as:

$$\boldsymbol{Q}_i = \boldsymbol{B}_i \cdot \boldsymbol{W}_Q, \boldsymbol{K}_i = \boldsymbol{B}_i \cdot \boldsymbol{W}_K, \boldsymbol{V}_i = \boldsymbol{B}_i \cdot \boldsymbol{W}_V, \tag{4}$$

where $W_Q \in \mathbb{R}^{k \times d_k}$, $W_K \in \mathbb{R}^{k \times d_k}$, and $W_V \in \mathbb{R}^{k \times d_v}$ are learned weight matrices; d_k is the dimensionality of the query and key spaces; d_v denotes the dimensionality of the value space. Then, we compute the attention weights by taking the dot product between the query and key matrices. To ensure numerical stability and manageable gradient magnitudes during training, the result is scaled by $\sqrt{d_k}$ and normalized using the softmax function. These attention weights are used to perform a A

270 weighted aggregation of the value matrix $V_i \in \mathbb{R}^{s \times d_v}$, which can be formulated as: 271

$$\operatorname{ttention}(\boldsymbol{Q}_i, \boldsymbol{K}_i, \boldsymbol{V}_i) = \operatorname{softmax}\left(\frac{\boldsymbol{Q}_i \boldsymbol{K}_i^T}{\sqrt{d_k}}\right) \boldsymbol{V}_i.$$
(5)

273 274

279

272

To comprehensively capture multiple facets of the feature subspace, we design multiple heads, each 275 with the same structure as described above. These heads generate different attention outputs from 276 different perspectives. The resulting attention outputs are then concatenated and passed through a 277 linear transformation to get $B'_i \in \mathbb{R}^{s \times d_{out}}$, which can be formulated as 278

$$\boldsymbol{B}_{i}^{\prime} = \operatorname{concat}(\operatorname{head}_{1}, \operatorname{head}_{2}, \dots, \operatorname{head}_{h})\boldsymbol{W}_{O}, \tag{6}$$

where h is the number of attention heads and $W_O \in \mathbb{R}^{(h \cdot d_v) \times d_{out}}$ is the output weight matrix. 281

After that, we concatenate B'_i and B_i to form a combined representation, which is then passed 282 283 through a fully connected layer to generate the prediction \hat{y}_i . This process can be formulated as follows: 284

$$\hat{\boldsymbol{y}}_i = FC(Concat(\boldsymbol{B}'_i; \boldsymbol{B}_i)). \tag{7}$$

(8)

This concatenation allows the evaluator to retain both original and context enhanced feature information for more effective prediction. When different feature subspaces are input into the evaluator, the same structure is used, and the weights are shared across all subspaces. This ensures consistency 288 and promotes generalization by learning common patterns. 289

290 291

285

287

4.3 **OPTIMIZATION**

292 **Pre-training.** A well-initialized contextual attention evaluator provides a strong foundation for eval-293 uation, allowing faster convergence and ensuring fair evaluation. To ensure an effective initialization for the contextual attention evaluator \mathcal{M} , we pre-train it using the original feature space \mathcal{F} as a foun-295 dational basis. Rather than employing the feature index optimizer and sample index optimizer, we 296 construct the feature subspaces $\mathcal{B}^{(0)}$ by randomly sampling the feature and sample indices. Then, we 297 subsequently input each feature subspace within $\mathcal{B}^{(0)}$ into the evaluator to perform the prediction. 298 The optimization objective is to minimize the discrepancy between the predicted and actual values, 299 which can be formulated as:

302

where $m{y}_i^{(0)}$ is the associated target label space of $m{B}_i^{(0)}; \hat{m{y}}_i^{(0)}$ is the predicted target label space; and s303 is number of feature space within $B_i^{(0)}$. After the model converges, the evaluator is initialized with 304 305 the parameters $\Theta_{\mathcal{M}}^{(0)}$. 306

 $\mathcal{L}_{ ext{intial}} = \sum_{i=1}^{s} \mathcal{L}_i(\boldsymbol{y}_i^{(0)}, \hat{\boldsymbol{y}}_i^{(0)}),$

Incremental Update. In the iterative feature optimization framework, consecutive iterations often 307 exhibit partial overlap in feature space information. This motivates us to incrementally update the 308 parameters of the contextual attention evaluator \mathcal{M} , enabling faster updates and accelerating the 309 entire feature space optimization process. 310

Specifically, in the t-th iteration, we begin by calculating the Fisher information to assess the 311 importance of parameters based on the previous learning iteration. Given the feature subspaces 312 $\mathcal{B}^{(t-1)} = \{ B_1^{(t-1)}, \cdots, B_q^{(t-1)} \}$, the Fisher information for the *j*-th parameter θ_j in the parameter 313 set $\Theta_{M}^{(t-1)}$ of the contextual attention evaluator is computed as: 314 315

$$\mathcal{G}(\boldsymbol{\theta}_{j}^{(t-1)}) = \frac{1}{s} \sum_{i=1}^{s} (\nabla_{\boldsymbol{\theta}_{j}} \log p(\boldsymbol{y}_{i}^{(t-1)} \mid \boldsymbol{B}_{i}^{(t-1)}; \boldsymbol{\Theta}_{\mathcal{M}}^{(t-1)})^{2},$$
(9)

316 317 318

where $\mathcal{G}(\boldsymbol{\theta}_{j}^{(t-1)})$ measures the importance of $\boldsymbol{\theta}_{j}$ based on its contribution to the evaluation task in the t-1 iteration. The term $\nabla_{\boldsymbol{\theta}_{j}} \log p(\boldsymbol{y}_{i}^{(t-1)} \mid \boldsymbol{B}_{i}^{(t-1)}; \boldsymbol{\Theta}_{\mathcal{M}}^{(t-1)})$ represents the logarithmic likelihood 319 320 321 gradient with respect to θ_j for conducting evaluation. A higher value of $\mathcal{G}(\theta_i^{(t-1)})$ indicates that 322 θ_i is crucial to perform an evaluation of $\mathcal{B}^{(t-1)}$ (Grosse & Martens, 2016). After that, we impose 323 constraints on parameter updates during the training of feature subspaces $\mathcal{B}^{(t)}$. The objective is to prevent forgetting shared evaluation knowledge during parameter updates while incorporating new evaluation insights. The final loss function in the *t*-th iteration is defined as:

$$\mathcal{L}_{\text{final}}(\boldsymbol{\Theta}_{\mathcal{M}}^{(t)}) = \mathcal{L}_{\mathcal{B}^{(t)}}(\boldsymbol{\Theta}_{\mathcal{M}}^{(t)}) + \frac{\lambda}{2} \sum_{j} \mathcal{G}(\boldsymbol{\theta}_{j}^{(t-1)}) \left(\boldsymbol{\theta}_{j}^{(t)} - \boldsymbol{\theta}_{j}^{(t-1)}\right)^{2}, \tag{10}$$

where $\mathcal{L}_{\mathcal{B}^{(t)}}(\Theta_{\mathcal{M}}^{(t)})$ represents the prediction loss for the current feature subspaces $\mathcal{B}^{(t)}$; $\theta_j^{(t)}$ and $\theta_j^{(t-1)}$ are the value of the *j*-th parameter from the parameter set of $\Theta_{\mathcal{M}}^{(t)}$ and $\Theta_{\mathcal{M}}^{(t-1)}$ respectively; λ is a regularization factor that balances incorporating new evaluation knowledge with preserving shared knowledge. During the optimization procedure, we minimize $\mathcal{L}_{\text{final}}(\Theta_{\mathcal{M}}^{(t)})$ to allow the contextual attention evaluator to efficiently capture the dynamics of the feature space, promoting faster convergence and more stable learning.

5 EXPERIMENTS

327 328

330

331 332

333

334

335 336

337 338

339

345

5.1 EXPERIMENTAL SETUP

Datasets. We conduct extensive experiments on 14 publicly available datasets from
UCI (Public, 2024b), OpenML (Public, 2024c) and Kaggle (Public, 2024a), consisting of 6 classification tasks and 6 regression tasks. A statistical overview of these
datasets is presented in Table 1. In this table, 'C' denotes dataset used for classification tasks and 'R' indicates datasets employed for represents regression tasks.

Evaluation Metrics. We use Mean Abso-346 lute Error (MAE), Root Mean Squared Error 347 (RMSE), and R-squared (R^2) to evaluate the 348 performance of regression tasks. Specifically, 349 $R^2 = 1 - \frac{\sum_{i=1}^{n} (y_i - \hat{y}_i)^2}{\sum_{i=1}^{n} (y_i - \bar{y})^2}$, where y_i and \hat{y}_i re-350 351 spectively represent the true label and predicted 352 label of x_i , and n is the number of samples. We 353 use Accuracy, Precision, Recall, and F1 score to evaluate the performance of classification tasks. 354

Baseline Algorithms. We apply EASE to two
iterative feature selection frameworks to validate its effectiveness and generalization capability: (i) RFE (Guyon et al., 2002) iteratively
eliminates the least important features from the
original set until a stopping criterion is met. (ii)
FLSR (Zhao et al., 2020) uses a single rein-

Table 1: Summary of the datasets.

Dataset	R/C	Samples	Features	Classes	Source
openml_607	R	1000	51	-	OpenML
openml_616	R	500	51	_	OpenML
openml_620	R	1000	26	_	OpenML
openml_586	R	1000	26	_	OpenML
airfoil	R	1503	6	_	OpenML
bike_share	R	10886	12	_	OpenML
wine_red	С	999	12	6	UCI
svmguide3	С	1243	22	2	OpenML
wine_white	С	4898	12	7	OpenML
spectf	С	267	45	2	UĈI
spam_base	С	4601	58	2	OpenML
mammography	С	11183	7	2	OpenML
spam_base	С	4601	58	2	OpenML
ÂmazonEA	С	32769	9	2	Kaggle
Nomao	С	34465	118	2	UCI

forced agent to perform feature selection with a restructured decision strategy. (iii) SDAE (Hassanieh & Chehade, 2024) is a state-of-the-art algorithm designed to select features used in unlabeled datasets without compromising information quality.

Additionally, we employ six widely-used ML algorithms as feature space evaluators during the it-365 erative optimization process to compare their experimental performance against EASE: (1) Linear 366 Regression / Logistic Regression (LR): Linear Regression (Su et al., 2012) models a linear relation-367 ship between the features and labels. Logistic Regression (Nusinovici et al., 2020) classifies data 368 by linearly combining input features and applying a logistic function to the result. **LR** refers to lin-369 ear regression for regression tasks and logistic regression for classification tasks. (2) Decision Tree 370 (DT): (Kim & Upneja, 2014) is a tree-like structure method, used to classify or predict data through 371 some rules. (3) Gradient Boosting Decision Tree (GBDT): (Li et al., 2023) builds an ensemble of de-372 cision trees sequentially to minimize errors. (4) Random Forest (RF): (Khajavi & Rastgoo, 2023) is 373 an ensemble learning method that constructs multiple decision trees. (5) Extreme Gradient Boosting 374 (XGB) (Asselman et al., 2023) combines the strengths of gradient boost with regularization tech-375 niques. In the testing phase, we use RF in all cases to report the performance of the refined feature space, as the model is stable and helps mitigate bias caused by the downstream model. For more ex-376 perimental details, we have provided hyperparameters and environmental settings in Appendix C.1 377 to improve the reproducibility of our work.

Dataset	R/C	Metrics	EASE	LR	DT	GBDT	RF	XGB
openml_607	R	$\begin{array}{c} \text{MAE} \downarrow \\ \text{RMSE} \downarrow \\ R^2 \uparrow \end{array}$	$\begin{array}{c} \textbf{0.271} \pm 0.028 \\ \textbf{0.344} \pm 0.036 \\ \textbf{0.873} \pm 0.025 \end{array}$	$\begin{array}{c} 0.342 \pm 0.058 \\ 0.439 \pm 0.070 \\ 0.805 \pm 0.087 \end{array}$	$\begin{array}{c} 0.290 \pm 0.020 \\ 0.374 \pm 0.019 \\ 0.847 \pm 0.033 \end{array}$	$\begin{array}{c} \underline{0.277} \pm 0.017 \\ 0.352 \pm 0.018 \\ \underline{0.863} \pm 0.010 \end{array}$	$\begin{array}{c} 0.301 \pm 0.025 \\ 0.398 \pm 0.041 \\ 0.835 \pm 0.026 \end{array}$	$ \begin{array}{r} 0.278 \pm \\ \underline{0.345} \pm \\ 0.863 \pm \\ \end{array} $
openml_616	R	$\begin{array}{c} \text{MAE} \downarrow \\ \text{RMSE} \downarrow \\ R^2 \uparrow \end{array}$	$\begin{array}{c} \textbf{0.302} \pm 0.039 \\ \textbf{0.389} \pm 0.048 \\ \textbf{0.840} \pm 0.035 \end{array}$	$\frac{\underline{0.315} \pm 0.035}{\underline{0.404} \pm 0.054} \\ \overline{0.833} \pm 0.028$	$\begin{array}{c} 0.365 \pm 0.011 \\ 0.469 \pm 0.023 \\ 0.799 \pm 0.039 \end{array}$	$\begin{array}{c} 0.321 \pm 0.044 \\ 0.418 \pm 0.063 \\ 0.826 \pm 0.044 \end{array}$	$\begin{array}{c} 0.367 \pm 0.037 \\ 0.466 \pm 0.038 \\ 0.791 \pm 0.039 \end{array}$	0.323 ± 0.406 ± <u>0.837</u> ±
openml_620	R	$\begin{array}{c} \text{MAE} \downarrow \\ \text{RMSE} \downarrow \\ R^2 \uparrow \end{array}$	$\begin{array}{c} \textbf{0.297} \pm 0.017 \\ \textbf{0.372} \pm 0.021 \\ \textbf{0.861} \pm 0.014 \end{array}$	$\begin{array}{c} 0.371 \pm 0.083 \\ 0.476 \pm 0.100 \\ 0.780 \pm 0.072 \end{array}$	$\begin{array}{c} 0.304 \pm 0.014 \\ 0.476 \pm 0.100 \\ 0.852 \pm 0.013 \end{array}$	$\begin{array}{c} 0.302 \pm 0.027 \\ \underline{0.380} \pm 0.033 \\ 0.848 \pm 0.031 \end{array}$	$\begin{array}{c} 0.313 \pm 0.015 \\ 0.395 \pm 0.016 \\ 0.846 \pm 0.007 \end{array}$	<u>0.298</u> 0.383 <u>0.855</u>
openml_586	R	$\begin{array}{c} \text{MAE} \downarrow \\ \text{RMSE} \downarrow \\ R^2 \uparrow \end{array}$	$\begin{array}{c} \textbf{0.277} \pm 0.011 \\ \textbf{0.357} \pm 0.019 \\ \textbf{0.875} \pm 0.015 \end{array}$	$\begin{array}{c} 0.313 \pm 0.057 \\ 0.405 \pm 0.074 \\ 0.818 \pm 0.062 \end{array}$	$\begin{array}{c} 0.294 \pm 0.024 \\ 0.368 \pm 0.033 \\ 0.862 \pm 0.033 \end{array}$	$\begin{array}{c} 0.291 \pm 0.024 \\ 0.373 \pm 0.033 \\ 0.861 \pm 0.022 \end{array}$	$\begin{array}{c} 0.283 \pm 0.020 \\ 0.369 \pm 0.025 \\ 0.862 \pm 0.022 \end{array}$	$\frac{0.281}{0.361}$ = $\frac{0.361}{0.872}$ =
mammograph	y C	Accuracy ↑ Precision ↑ F1 ↑ Recall ↑	$\begin{array}{c} \textbf{0.989} \pm 0.003 \\ 0.947 \pm 0.042 \\ \textbf{0.818} \pm 0.041 \\ \textbf{0.831} \pm 0.043 \end{array}$	$\begin{array}{c} 0.978 \pm 0.004 \\ 0.837 \pm 0.036 \\ 0.653 \pm 0.035 \\ 0.603 \pm 0.029 \end{array}$	$\begin{array}{c} 0.984 \pm 0.004 \\ \textbf{0.976} \pm 0.021 \\ 0.713 \pm 0.095 \\ 0.651 \pm 0.078 \end{array}$	$\begin{array}{c} \underline{0.987} \pm 0.002 \\ 0.949 \pm 0.018 \\ \underline{0.803} \pm 0.041 \\ \underline{0.733} \pm 0.045 \end{array}$	$\begin{array}{c} 0.985 \pm 0.005 \\ 0.947 \pm 0.010 \\ 0.757 \pm 0.034 \\ 0.685 \pm 0.035 \end{array}$	0.985 = 0.953 = 0.736 = 0.662 =
spectf	С	Accuracy ↑ Precision ↑ F1 ↑ Recall ↑	$\begin{array}{c} \textbf{0.825} \pm 0.041 \\ \underline{0.631} \pm 0.248 \\ \textbf{0.543} \pm 0.054 \\ 0.515 \pm 0.092 \end{array}$	$\begin{array}{c} 0.737 \pm 0.028 \\ 0.467 \pm 0.209 \\ 0.454 \pm 0.068 \\ 0.518 \pm 0.036 \end{array}$	$\begin{array}{c} 0.781 \pm 0.049 \\ \textbf{0.643} \pm 0.237 \\ \underline{0.512} \pm 0.077 \\ \textbf{0.542} \pm 0.038 \end{array}$	$\begin{array}{c} 0.795 \pm 0.055 \\ 0.466 \pm 0.147 \\ 0.478 \pm 0.079 \\ \underline{0.522} \pm 0.044 \end{array}$	$\begin{array}{c} \underline{0.815} \pm 0.057 \\ \hline 0.459 \pm 0.117 \\ 0.473 \pm 0.058 \\ 0.514 \pm 0.028 \end{array}$	0.766 ± 0.482 ± 0.450 ± 0.509 ±
AmazonEA	С	Accuracy ↑ Precision ↑ F1 ↑ Recall ↑	$\begin{array}{c} \textbf{0.963} \pm 0.003 \\ \underline{0.782} \pm 0.199 \\ \textbf{0.500} \pm 0.001 \\ \textbf{0.503} \pm 0.001 \end{array}$	$\begin{array}{c} 0.941 \pm 0.004 \\ 0.670 \pm 0.246 \\ 0.488 \pm 0.005 \\ 0.501 \pm 0.002 \end{array}$	$\begin{array}{c} \underline{0.944} \pm 0.002 \\ \hline 0.772 \pm 0.245 \\ 0.489 \pm 0.003 \\ 0.502 \pm 0.002 \end{array}$	$\begin{array}{c} 0.943 \pm 0.002 \\ 0.489 \pm 0.001 \\ \underline{0.490} \pm 0.002 \\ \underline{0.502} \pm 0.001 \end{array}$	$\begin{array}{c} 0.944 \pm 0.003 \\ \textbf{0.872} \pm 0.199 \\ 0.489 \pm 0.003 \\ \underline{0.502} \pm 0.001 \end{array}$	0.942 = 0.611 = 0.488 = 0.501 =
Nomao	С	Accuracy ↑ Precision ↑ F1 ↑ Recall ↑	$0.952 \pm 0.005 \\ 0.947 \pm 0.004 \\ 0.936 \pm 0.007 \\ 0.922 \pm 0.006$	$\begin{array}{c} 0.942 \pm 0.002 \\ 0.940 \pm 0.004 \\ 0.927 \pm 0.003 \\ 0.916 \pm 0.003 \end{array}$	$\frac{0.944}{0.941} \pm 0.004 \\ \frac{0.941}{0.930} \pm 0.003 \\ \frac{0.930}{0.920} \pm 0.005 \\ 0.920 \pm 0.006 \\ 0.900 $	$\begin{array}{c} 0.940 \pm 0.003 \\ 0.937 \pm 0.002 \\ 0.924 \pm 0.004 \\ 0.914 \pm 0.006 \end{array}$	$\begin{array}{c} 0.941 \pm 0.003 \\ 0.936 \pm 0.003 \\ 0.926 \pm 0.004 \\ 0.917 \pm 0.005 \end{array}$	0.937 = 0.935 = 0.922 = 0.911 =

Table 2: Overall performance comparison. The best results are highlighted in **bold**, and the secondbest results are <u>underlined</u>. (\uparrow indicates that a higher value of the metric corresponds to better performance, while \downarrow denotes that a lower value of the metric indicates better performance.)



Figure 3: Time complexity comparison of different feature space evaluators across various datasets.

5.2 EXPERIMENTAL RESULTS

404

413

418 419

420

5.2.1 OVERALL COMPARISON

421 This experiment aims to answer: Can EASE accurately assess feature space quality to produce an 422 *effective feature space*? We choose RFE for iterative FS and adopt EASE as the feature space evalu-423 ator. To compare the performance difference, we replace the evaluator with LR, DT, GBDT, RF and 424 XGB respectively. We report the testing performance of the refined feature space using RF. Table 2 425 and Appendix C.2 shows the comparison results in terms of different evaluation metrics according 426 to the task type. We can find that EASE outperforms other baselines in most cases. For classi-427 fication, EASE can improve by approximately 3% compared to other baselines. For regression, 428 EASE demonstrates the most superior performance. The underlying driver is that our information 429 decoupling strategy and context-aware evaluator, which allow the evaluator to focus on the most challenging aspects of the feature space. This results in a fairer evaluation, leading to a more effec-430 tive refinement strategy and ultimately producing a more optimized feature space. In summary, this 431 experiment shows that EASE effectively evaluates feature space quality for better feature spaces.



Figure 5: Comparison of refinement performance of feature space evaluators within FLSR.

5.2.2 EFFICIENCY COMPARISON

This experiment aims to answer: *Is* EASE *more efficient compared to other feature space evalua tors?* We compare the training time of EASE with other feature space evaluators, including GBDT, RF, and XGB. Figure 3 shows the comparison results in terms of cumulative time. An interesting observation is that using EASE can significantly reduce the cumulative time costs compared to other baselines. A potential reason is that our incremental parameter update strategy enables the feature space evaluator to quickly capture evolving patterns in the feature space, thereby accelerating the feature optimization process. Additional experiments comparing time complexity across different datasets are provided in Figure C.3 in the Appendix. To sum up, this experiment demonstrates that EASE can efficiently assess feature space quality, thanks to its adaptive parameter update strategy.

5.2.3 THE EFFECTIVENESS OF EASE FOR FEATURE SPACE REFINEMENT

This experiment aims to answer: Is the quality of the feature space refined by EASE superior to the original feature space? We compare the prediction performance between the original feature space and the space refined by EASE using various downstream predictors, including LR, DT, GBDT, and RF. Figure 4 shows the comparison results in terms of Accuracy and MAE according to the task type. We find that the feature space produced by EASE outperforms the original in most cases across various predictors. In particular, the refinement by EASE outperforms the original feature space by 20% on the spectf dataset. This observation suggests that incorporating EASE into the iterative feature space optimization framework provides effective guidance for obtaining a better feature space. Additionally, the contextual attention evaluator comprehensively captures the intrinsic traits of the feature space, leading to robust performance across various datasets and predictors.

5.2.4 EASE'S PERFORMANCE IN DIFFERENT ITERATIVE FRAMEWORKS

This experiment aims to answer: Is EASE generalizable and applicable across different iterative feature space optimization algorithms? We apply EASE to a reinforced feature selection framework FLSR. To develop controlling groups, we replace the feature space evaluator within FLSR with LR, DT, GBDT, and RF, respectively. Figure 5 shows the comparison results across different datasets in terms of accuracy and performance standard deviation. We observe that EASE beats other baselines across various datasets. Especially, EASE improve the accuracy more than 5% in all situations. This observation highlights the strong generalizability and applicability of EASE. The underlying driver is that the feature index optimizer of EASE offers flexibility, allowing it to adapt to various iterative feature optimization frameworks. In summary, this experiment demonstrates that EASE exhibits

486	Table 3: Comparison of different EASE variants in terms of Accuracy or MAE. The best results
487	are highlighted in bold , and the second-best results are <u>underlined</u> . (↑ indicates that a higher value
488	of the metric corresponds to better performance, while \downarrow denotes that a lower value of the metric
489	indicates better performance.)

Dataset	R/C	Metric	EASE	$EASE^{-FC}$	$EASE^{-IT}$	$EASE^{-PT}$
openml_607	R	MAE↓	$\textbf{0.271} \pm 0.028$	$\underline{0.289} \pm 0.028$	0.343 ± 0.021	0.313 ± 0.035
openml_616	R	MAE 🗸	$\textbf{0.302} \pm 0.039$	0.339 ± 0.035	0.435 ± 0.038	0.561 ± 0.082
openml_620	R	MAE 🗸	$\textbf{0.297} \pm 0.017$	$\overline{0.520}\pm0.018$	0.383 ± 0.018	$\underline{0.318} \pm 0.017$
openml_586	R	MAE 🗸	$\textbf{0.277} \pm 0.011$	0.348 ± 0.016	0.276 ± 0.014	0.364 ± 0.005
airfoil	R	MAE 🗸	$\textbf{0.337} \pm 0.009$	0.751 ± 0.029	0.712 ± 0.028	$\underline{0.362} \pm 0.004$
bike_share	R	MAE \downarrow	$\textbf{0.033} \pm 0.001$	0.035 ± 0.002	$\underline{0.034} \pm 0.001$	$\underline{0.034} \pm 0.001$
wine_red	С	Accuracy ↑	$\textbf{0.637} \pm 0.018$	$\underline{0.589} \pm 0.055$	0.539 ± 0.043	0.539 ± 0.028
svmguide3	С	Accuracy	$\textbf{0.846} \pm 0.034$	0.817 ± 0.032	0.828 ± 0.014	$\underline{0.833} \pm 0.010$
wine_white	С	Accuracy ↑	$\textbf{0.578} \pm 0.017$	0.557 ± 0.019	0.576 ± 0.014	0.552 ± 0.011
spam_base	С	Accuracy ↑	$\textbf{0.939} \pm 0.008$	0.918 ± 0.016	0.928 ± 0.007	0.931 ± 0.003
mammography	С	Accuracy ↑	$\textbf{0.989} \pm 0.003$	0.981 ± 0.004	0.985 ± 0.004	$\underline{0.986} \pm 0.003$
spectf	С	Accuracy ↑	$\textbf{0.825} \pm 0.041$	$\underline{0.815}\pm0.065$	0.756 ± 0.064	0.781 ± 0.041

strong adaptability to iterative feature space optimization frameworks and excellent generalizability across different optimization algorithms. We Additionally test the effectiveness of EASE and applied it to the state-of-the-art FS algorithm. The detailed results are provided in Appendix C.9.

5.2.5 THE IMPACT OF EACH TECHNICAL COMPONENT

508 This experiment aims to answer: How does each technical component in EASE impact its perfor-509 *mance*? We investigate the effects of pre-training, incremental training, and feature-sample subspace construction in EASE. We develop $EASE^{-PT}$, $EASE^{-IT}$, and $EASE^{-FC}$ by removing the 510 511 pre-training, incremental training, and feature-sample subspace construction steps from EASE re-512 spectively. Table 3 shows the comparison results among different EASE variants. We observe that 513 EASE outperforms $EASE^{-PT}$, highlighting the importance of the pre-training step in providing 514 a strong foundation for objective evaluation. Additionally, EASE surpasses $EASE^{-IT}$, indicating 515 that the incremental parameter updating mechanism effectively captures the evolving patterns of 516 feature space optimization, leading to an improved feature space. Moreover, EASE outperforms EASE-FC, showing that the information decoupling strategy reduces comprehension complexity, 517 allowing for better capture of feature space characteristics and leading to improved evaluation and 518 feature space quality. We also compare the time complexity of each component in Appendix C.6. In 519 conclusion, this experiment reflects that each component in EASE is indispensable and significant. 520

For additional experiments and case studies, please refer to the Appendix, which further demonstratethe superiority, efficiency, and generalization capability of EASE.

523 524

502 503

504

505 506 507

6 CONCLUSION

525 526

In this paper, we propose a generalized adaptive feature space evaluator EASE for iterative feature 527 space optimization. EASE consists of two key components: feature-sample subspace generator and 528 contextual attention evaluator. The first component decouples the complex information within the 529 feature space to generate diverse feature subspaces by the cooperation of the feature index optimizer 530 and sample index optimizer. This enhances the ability of the subsequent evaluator to capture the most 531 challenging information for more accurate feature space evaluation. The second component captures 532 the intrinsic complexity of feature-sample interactions using a weight-sharing contextual-attention 533 evaluator to ensure fair and accurate evaluation. Considering the information overlap across consec-534 utive iterations, we incrementally update the evaluator's parameters to retain past knowledge while 535 incorporating new insights. This allows the evaluator to efficiently capture the evolving patterns 536 of the feature space. Extensive experimental results have demonstrated that EASE has achieved 537 superior performance compared to other baselines. In addition, EASE exhibits strong adaptability, generalization, and robustness in various iterative feature optimization frameworks. In the future, 538 we will focus on further enhancing the generalization capability of EASE to enable it to effectively handle distribution shifts and perform robustly across different types of datasets.

540 REFERENCES

548

565

571

- Rahaf Aljundi, Punarjay Chakravarty, and Tinne Tuytelaars. Expert gate: Lifelong learning with
 a network of experts. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 3366–3375, 2017.
- Nisha Arora and Pankaj Deep Kaur. A bolasso based consistent feature selection enabled random
 forest classification algorithm: An application to credit risk assessment. *Applied Soft Computing*,
 86:105936, 2020.
- Amal Asselman, Mohamed Khaldi, and Souhaib Aammou. Enhancing the prediction of student performance based on the machine learning xgboost algorithm. *Interactive Learning Environments*, 31(6):3360–3379, 2023.
- Xiyang Dai, Yinpeng Chen, Bin Xiao, Dongdong Chen, Mengchen Liu, Lu Yuan, and Lei Zhang.
 Dynamic head: Unifying object detection heads with attentions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 7373–7382, June 2021.
- Burcu F Darst, Kristen C Malecki, and Corinne D Engelman. Using recursive feature elimination in random forest to account for correlated variables in high dimensional data. *BMC genetics*, 19: 1–6, 2018.
- Nicholas Sean Escanilla, Lisa Hellerstein, Ross Kleiman, Zhaobin Kuang, James Shull, and David
 Page. Recursive feature elimination by sensitivity testing. In 2018 17th IEEE International
 Conference on Machine Learning and Applications (ICMLA), pp. 40–47, 2018. doi: 10.1109/
 ICMLA.2018.00014.
- Roger Grosse and James Martens. A kronecker-factored approximate fisher matrix for convolution
 layers. In *International Conference on Machine Learning*, pp. 573–582. PMLR, 2016.
- Isabelle Guyon, Jason Weston, Stephen Barnhill, and Vladimir Vapnik. Gene selection for cancer
 classification using support vector machines. *Machine learning*, 46:389–422, 2002.
- Wael Hassanieh and Abdallah Chehade. Selective deep autoencoder for unsupervised feature selection. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pp. 12322–12330, 2024.
- Htet Htet Htun, Michael Biehl, and Nicolai Petkov. Survey of feature selection and extraction techniques for stock market prediction. *Financial Innovation*, 9(1):26, 2023.
- David Isele and Akansel Cosgun. Selective experience replay for lifelong learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32, 2018.
- Weikuan Jia, Meili Sun, Jian Lian, and Sujuan Hou. Feature dimensionality reduction: a review. *Complex & Intelligent Systems*, 8(3):2663–2693, 2022.
- Hamed Khajavi and Amir Rastgoo. Predicting the carbon dioxide emission caused by road transport using a random forest (rf) model combined by meta-heuristic algorithms. *Sustainable Cities and Society*, 93:104503, 2023.
- Soo Y Kim and Arun Upneja. Predicting restaurant financial distress using decision tree and adaboosted decision tree models. *Economic Modelling*, 36:354–362, 2014.
- James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A
 Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al. Overcom ing catastrophic forgetting in neural networks. *Proceedings of the national academy of sciences*, 114(13):3521–3526, 2017.
- Zhen Li, Tieding Lu, Kegen Yu, and Jie Wang. Interpolation of gnss position time series using gbdt, xgboost, and rf machine learning algorithms and models error analysis. *Remote Sensing*, 15(18): 4374, 2023.
- ⁵⁹³ Zhizhong Li and Derek Hoiem. Learning without forgetting. *IEEE transactions on pattern analysis and machine intelligence*, 40(12):2935–2947, 2017.

ing systems, 32(6):2306–2319, 2021b.

600

603

634

635

636

637

594 595 596	Kunpeng Liu, Yanjie Fu, Le Wu, Xiaolin Li, Charu Aggarwal, and Hui Xiong. Automated feature selection: A reinforcement learning perspective. <i>IEEE Transactions on Knowledge and Data Engineering</i> , 35(3):2272–2284, 2021a.
597	L Liu 7 Kuong V Chan Ib Vug W Vang and W 7bang Indat. In defense of electic weight
598 599	consolidation for incremental object detection. <i>IEEE transactions on neural networks and learn</i> -

- Zhaogeng Liu, Jielong Yang, Li Wang, and Yi Chang. A novel relation aware wrapper method for
 feature selection. *Pattern Recognition*, 140:109566, 2023.
- Kaouther Messaoud, Nachiket Deo, Mohan M Trivedi, and Fawzi Nashashibi. Trajectory prediction
 for autonomous driving based on multi-head attention with joint agent-map representation. In
 2021 IEEE Intelligent Vehicles Symposium (IV), pp. 165–170. IEEE, 2021.
- Takafumi Moriya, Ryo Masumura, Taichi Asami, Yusuke Shinohara, Marc Delcroix, Yoshikazu
 Yamaguchi, and Yushi Aono. Progressive neural network-based knowledge transfer in acoustic
 models. In 2018 Asia-Pacific Signal and Information Processing Association Annual Summit and
 Conference (APSIPA ASC), pp. 998–1002. IEEE, 2018.
- Bach Hoai Nguyen, Bing Xue, and Mengjie Zhang. A survey on swarm intelligence approaches to feature selection in data mining. *Swarm and Evolutionary Computation*, 54:100663, 2020.
- Babak Nouri-Moghaddam, Mehdi Ghazanfari, and Mohammad Fathian. A novel multi-objective
 forest optimization algorithm for wrapper feature selection. *Expert Systems with Applications*, 175:114737, 2021.
- Simon Nusinovici, Yih Chung Tham, Marco Yu Chak Yan, Daniel Shu Wei Ting, Jialiang Li, Charumathi Sabanayagam, Tien Yin Wong, and Ching-Yu Cheng. Logistic regression was as good as machine learning for predicting major chronic diseases. *Journal of clinical epidemiology*, 122: 56–69, 2020.
- Arif Mudi Priyatno, Triyanna Widiyaningtyas, et al. A systematic literature review: Recursive feature elimination algorithms. *JITK (Jurnal Ilmu Pengetahuan dan Teknologi Komputer)*, 9(2): 196–207, 2024.
- Public. Kaggle dataset download. [EB/OL]. https://www.kaggle.com/c/
 amazon-employee-access-challenge/data, 2024a.
- 628 Public. Openml dataset download. [EB/OL]. https://www.openml.org, 2024b.
- Public. Uci dataset download. [EB/OL]. https://archive.ics.uci.edu/, 2024c.
- Nicholas Pudjihartono, Tayaza Fadason, Andreas W Kempa-Liehr, and Justin M O'Sullivan. A
 review of feature selection methods for machine learning-based disease risk prediction. *Frontiers in Bioinformatics*, 2:927312, 2022.
 - Jathushan Rajasegaran, Munawar Hayat, Salman H Khan, Fahad Shahbaz Khan, and Ling Shao. Random path selection for continual learning. *Advances in neural information processing systems*, 32, 2019.
- Jesse Read, Albert Bifet, Bernhard Pfahringer, and Geoff Holmes. Batch-incremental versus instance-incremental learning in dynamic and evolving data. In *Advances in Intelligent Data Analysis XI: 11th International Symposium, IDA 2012, Helsinki, Finland, October 25-27, 2012. Proceedings 11*, pp. 313–323. Springer, 2012.
- Joan Serra, Didac Suris, Marius Miron, and Alexandros Karatzoglou. Overcoming catastrophic
 forgetting with hard attention to the task. In *International conference on machine learning*, pp. 4548–4557. PMLR, 2018.
- Jeng-Lun Shieh, Qazi Mazhar ul Haq, Muhamad Amirul Haq, Said Karam, Peter Chondro, De-Qin
 Gao, and Shanq-Jang Ruan. Continual learning strategy in one-stage object detection framework
 based on experience replay for autonomous driving vehicle. *Sensors*, 20(23):6777, 2020.

- 648 Dongsub Shim, Zheda Mai, Jihwan Jeong, Scott Sanner, Hyunwoo Kim, and Jongseong Jang. On-649 line class-incremental continual learning with adversarial shapley value. In Proceedings of the 650 AAAI Conference on Artificial Intelligence, volume 35, pp. 9630–9638, 2021. 651 Hanul Shin, Jung Kwon Lee, Jaehong Kim, and Jiwon Kim. Continual learning with deep generative 652 replay. Advances in neural information processing systems, 30, 2017. 653 654 Xiaogang Su, Xin Yan, and Chih-Ling Tsai. Linear regression. Wiley Interdisciplinary Reviews: 655 Computational Statistics, 4(3):275–294, 2012. 656 Zewei Sun, Shujian Huang, Hao-Ran Wei, Xin-yu Dai, and Jiajun Chen. Generating diverse trans-657 lation by manipulating multi-head attention. In Proceedings of the AAAI Conference on Artificial 658 Intelligence, volume 34, pp. 8976–8983, 2020. 659 Anastasios Temenos, Nikos Temenos, Maria Kaselimi, Anastasios Doulamis, and Nikolaos 661 Doulamis. Interpretable deep learning framework for land use and land cover classification in 662 remote sensing using shap. *IEEE Geoscience and Remote Sensing Letters*, 20:1–5, 2023. 663 A Vaswani. Attention is all you need. Advances in Neural Information Processing Systems, 2017. 664 665 Amukta Malyada Vommi and Tirumala Krishna Battula. A hybrid filter-wrapper feature selection 666 using fuzzy knn based on bonferroni mean for medical datasets classification: A covid-19 case 667 study. Expert Systems with Applications, 218:119612, 2023. 668 Dongjie Wang, Yanjie Fu, Kunpeng Liu, Xiaolin Li, and Yan Solihin. Group-wise reinforcement fea-669 ture generation for optimal and explainable representation space reconstruction. In Proceedings of 670 the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, pp. 1826–1834, 671 2022. 672 Xubin Wang, Haojiong Shangguan, Fengyi Huang, Shangrui Wu, and Weijia Jia. Mel: efficient 673 multi-task evolutionary learning for high-dimensional feature selection. *IEEE Transactions on* 674 Knowledge and Data Engineering, 2024. 675 676 Meng Xiao, Dongjie Wang, Min Wu, Kunpeng Liu, Hui Xiong, Yuanchun Zhou, and Yanjie Fu. 677 Traceable group-wise self-optimizing feature transformation learning: A dual optimization per-678 spective. ACM Transactions on Knowledge Discovery from Data, 18(4):1–22, 2024. 679 Rizgar Zebari, Adnan Abdulazeez, Diyar Zeebaree, Dilovan Zebari, and Jwan Saeed. A compre-680 hensive review of dimensionality reduction techniques for feature selection and feature extraction. 681 Journal of Applied Science and Technology Trends, 1(1):56–70, 2020. 682 683 Xiaosa Zhao, Kunpeng Liu, Wei Fan, Lu Jiang, Xiaowei Zhao, Minghao Yin, and Yanjie Fu. Simpli-684 fying reinforced feature selection via restructured choice strategy of single agent. In 2020 IEEE 685 International conference on data mining (ICDM), pp. 871–880. IEEE, 2020. 686 Fei Zhu, Zhen Cheng, Xu-Yao Zhang, and Cheng-lin Liu. Class-incremental learning via dual 687 augmentation. Advances in Neural Information Processing Systems, 34:14306–14318, 2021. 688 689 Yongbin Zhu, Wenshan Li, and Tao Li. A hybrid artificial immune optimization for high-690 dimensional feature selection. Knowledge-Based Systems, 260:110111, 2023. 691 692 693 694 696 697 699
- 700

A TABLE OF NOTATIONS

We provide key notations in Table 4 to enhance the comprehension of the EASE methodology. Among these notations, t represents the t-th optimization of the feature space algorithm.

	Table 4: Notations for EASE.							
Notations	Interpretation							
$\mathcal{M}^{(t)}$	Model of EASE.							
$oldsymbol{\Theta}_{\mathcal{M}}^{(t)}$	Parameter of EASE.							
$oldsymbol{\mathcal{F}}^{(t)}$	Feature space.							
$\mathcal{L}(\mathcal{M}(\cdot), oldsymbol{y})$	Prediction loss.							
$oldsymbol{f}^{(t)}$	Feature index subset.							
$\mathbb{I}^{(t)}$	Sample index set.							
$\mathcal{B}^{(t)}$	Feature subspaces.							
$\mathcal{G}_t(oldsymbol{ heta}_j)$,	Fisher information for parameter θ_j .							
$\mathcal{L}_{ ext{final}}(oldsymbol{\Theta}_{\mathcal{M}}^{(t)})$	The final loss includes the incremental loss and prediction loss.							

B ALGORITHMS OF EASE

Algorithm 1 is the pseudo-code for feature subspace construction based on the feature space $\mathcal{F}^{(t)}$ and the sample loss $\mathcal{L}_i^{(t-1)}$ from the previous iteration. Specifically, we first obtain an optimized feature index subset $f^{(t)}$ by Feature Index Optimizer. Then, a sampling probability distribution p_i and CDF C_i are constructed. Next, n samples are sampled to obtain the sample index set $\mathbb{I}^{(t)}$. Finally, the feature subspace $\mathcal{B}^{(t)}$ and $y^{(t)}$ is constructed based on the $\mathbb{I}^{(t)}$, s and $f^{(t)}$.

Alg	orithm 1 Feature Subspace Construction Algorithm.
1:	Input: feature space $\mathcal{F}^{(t)}$, fixed batch size <i>s</i> , sample loss $\mathcal{L}_i^{(t-1)}$.
2:	Output: $\mathcal{B}^{(t)}, y^{(t)}$.
3:	$f^{(t)} \leftarrow$ Feature index subset.
4:	$p_i \leftarrow$ Sampling probability distribution.
5:	$C_i \leftarrow \text{CDF}.$
6:	for $i = 0$ to n do
7:	generate a random number $r \in [0, 1]$,
8:	find the first sample x_i that $C_i \ge r$,
9:	add the sample index i to the set $\mathbb{I}^{(t)}$.
10:	end for
11:	$\mathcal{B}^{(t)}, \boldsymbol{y}^{(t)} \leftarrow \text{Base on } \boldsymbol{f}^{(t)}, s \text{ and } \mathbb{I}^{(t)}.$

Algorithm 2 is the pseudo-code for obtaining a robust evaluator. Specifically, we first construct a Contextual Attention Evaluator \mathcal{M} and initialize its parameters Θ . Then, we use Algorithm 1 to obtain the feature subspace $\mathcal{B}^{(0)}$ of all features and to pre-train the \mathcal{M} . For the *t*-th evaluation, we obtain the feature subspace $\mathcal{B}^{(t)}$ at first, and then estimate the importance of each parameter $\mathcal{G}_{t-1}(\theta_j)$, reducing the update magnitude for important parameters in previous iteration. Finally, we formulate the our objective function $\mathcal{L}_{\text{final}}(\Theta_{\mathcal{M}}^{(t)})$ and minimize it.

756 Algorithm 2 Generalized Adaptive Feature Space Evaluator Algorithm. 1: Input: Dataset $\mathbb{D} = \langle \mathcal{F}, y \rangle$, optimization iterations T, fixed batch size s 758 2: **Output:** Evaluator \mathcal{M} . 759 3: $\mathcal{M} \leftarrow \text{Contextual attention evaluator.}$ 760 4: $\Theta \leftarrow$ Initialize the parameters. 761 $\mathcal{B}^{(0)}, \boldsymbol{y}^{(0)} \leftarrow$ Feature subspace construction. 5: 762 $\boldsymbol{\Theta}_{\mathcal{M}}^{(0)} \leftarrow \text{Pre-training based on } \mathcal{B}^{(0)}.$ 6: 763 7: for t = 0 to T do 764 $\mathcal{B}^{(t)}, \boldsymbol{y}^{(t)} \leftarrow$ feature subspace construction, 8: 765 9: $\mathcal{G}_{t-1}(\boldsymbol{\theta}_j) \leftarrow \text{estimate parameter importance,}$ 766 $\mathcal{L}_{\text{final}}(\mathbf{\Theta}_{\mathcal{M}}^{(t)}) \leftarrow \text{final loss}$. 10: 767 $\boldsymbol{\Theta}_{\mathcal{M}}^{(t)} \leftarrow \text{train } \mathcal{M} \text{ by minimizing } \mathcal{L}_{\text{final}}(\boldsymbol{\Theta}_{\mathcal{M}}^{(t)}).$ 11: 768 12: end for 769

C ADDITIONAL EXPERIMENT RESULTS

C.1 EXPERIMENTAL SETUP

770 771

772 773

774

780 781

785

786

787 788 789

790 791

793

794

802

803 804

805 806 807

808

809

775Hyperparameters, Source Code and Reproducibility. We limit pre-training and incremental
training to 50 and 200, respectively. We employ an early stopping strategy, stopping the training
process when the loss does not decrease for 10 consecutive epochs. In all experiments, we use the
Adam optimizer and a learning rate decay strategy to accelerate the convergence. Specifically, the
learning rate for the t'-th training iteration is:

$$l(t') = l(t'_0) \times p^{\left\lfloor \frac{t'}{u} \right\rfloor},\tag{11}$$

where l(t') and $l(t'_0)$ is current and initial learning rate, p is the decay factor applied every u iterations, and $\lfloor \cdot \rfloor$ is floor operation. And we set $l(t'_0) = 0.001$, p = 0.9 and u = 30. All experiments run 10 times and calculate the value of mean and standard deviation.

Environmental Settings All experiments were conducted on the macOS Sonoma 14.0 operating system, Apple M3 Chip with 8 cores (4 performance and 4 efficiency), and 8GB of RAM, with the framework of Python 3.8.19 and TensorFlow 2.13.0.



Figure 6: Time complexity comparison of different feature space evaluators across various datasets.

C.2 ADDITIONAL OVERALL COMPARISON

We additionally compare overall comparison across all datasets. We choose RFE for iterative FS and adopt EASE as the feature space evaluator. To compare the performance difference, we replace the

Dataset	R/0	C Metrics	EASE	LR	DT	GBDT	RF	XGB
airfoil	R	$\begin{array}{c} \text{MAE} \downarrow \\ \text{RMSE} \downarrow \\ R^2 \uparrow \end{array}$	$\begin{array}{c} \textbf{0.337} \pm 0.010 \\ \textbf{0.440} \pm 0.017 \\ \textbf{0.807} \pm 0.010 \end{array}$	$\begin{array}{c} 0.348 \pm 0.011 \\ 0.449 \pm 0.014 \\ 0.801 \pm 0.012 \end{array}$	$\begin{array}{c} 0.364 \pm 0.014 \\ 0.475 \pm 0.028 \\ 0.767 \pm 0.012 \end{array}$	$\begin{array}{c} \underline{0.339} \pm 0.014 \\ \underline{0.448} \pm 0.016 \\ 0.799 \pm 0.026 \end{array}$	$\begin{array}{c} 0.347 \pm 0.017 \\ 0.454 \pm 0.014 \\ 0.783 \pm 0.013 \end{array}$	0.347 0.475 <u>0.802</u>
bike_share	R	$\begin{array}{c} \text{MAE} \downarrow \\ \text{RMSE} \downarrow \\ R^2 \uparrow \end{array}$	$\begin{array}{c} \textbf{0.033} \pm 0.001 \\ \textbf{0.051} \pm 0.002 \\ \textbf{0.998} \pm 0.000 \end{array}$	$\begin{array}{c} \underline{0.033} \pm 0.002 \\ 0.053 \pm 0.002 \\ \underline{0.997} \pm 0.000 \end{array}$	$\begin{array}{c} 0.034 \pm 0.002 \\ 0.053 \pm 0.003 \\ 0.997 \pm 0.001 \end{array}$	$\begin{array}{c} 0.036 \pm 0.000 \\ 0.057 \pm 0.001 \\ 0.997 \pm 0.001 \end{array}$	$\begin{array}{c} 0.034 \pm 0.001 \\ \underline{0.052} \pm 0.001 \\ 0.997 \pm 0.001 \end{array}$	0.037 0.052 0.996
wine_red	С	Accuracy ↑ Precision↑ F1 ↑ Recall ↑	$\begin{array}{c} \textbf{0.637} \pm 0.018 \\ \textbf{0.386} \pm 0.062 \\ \textbf{0.332} \pm 0.057 \\ \textbf{0.334} \pm 0.060 \end{array}$	$\begin{array}{c} \underline{0.617} \pm 0.042 \\ \underline{0.311} \pm 0.099 \\ \overline{0.260} \pm 0.041 \\ 0.276 \pm 0.038 \end{array}$	$\begin{array}{c} 0.599 \pm 0.043 \\ 0.296 \pm 0.018 \\ 0.290 \pm 0.017 \\ 0.294 \pm 0.015 \end{array}$	$\begin{array}{c} 0.596 \pm 0.026 \\ 0.264 \pm 0.036 \\ 0.263 \pm 0.027 \\ 0.276 \pm 0.025 \end{array}$	$\begin{array}{c} 0.613 \pm 0.020 \\ 0.307 \pm 0.066 \\ \underline{0.300} \pm 0.045 \\ \underline{0.308} \pm 0.031 \end{array}$	0.588 0.285 0.269 0.283
svmguide3	С	Accuracy ↑ Precision ↑ F1 ↑ Recall ↑	$\begin{array}{c} \textbf{0.846} \pm 0.034 \\ \textbf{0.852} \pm 0.052 \\ 0.646 \pm 0.026 \\ \textbf{0.679} \pm 0.036 \end{array}$	$\begin{array}{c} 0.815 \pm 0.015 \\ 0.805 \pm 0.033 \\ 0.655 \pm 0.048 \\ 0.635 \pm 0.038 \end{array}$	$\begin{array}{c} 0.816 \pm 0.020 \\ 0.823 \pm 0.034 \\ 0.673 \pm 0.040 \\ 0.651 \pm 0.033 \end{array}$	$\begin{array}{c} 0.813 \pm 0.024 \\ 0.787 \pm 0.044 \\ \underline{0.677} \pm 0.059 \\ 0.655 \pm 0.047 \end{array}$	$\begin{array}{c} 0.817 \pm 0.042 \\ 0.823 \pm 0.044 \\ \textbf{0.691} \pm 0.075 \\ \underline{0.670} \pm 0.060 \end{array}$	0.818 0.827 0.656 0.637
wine_white	С	Accuracy ↑ Precision ↑ F1 ↑ Recall ↑	$\begin{array}{c} \textbf{0.573} \pm 0.017 \\ \textbf{0.307} \pm 0.074 \\ 0.238 \pm 0.016 \\ \textbf{0.248} \pm 0.018 \end{array}$	$\begin{array}{c} 0.531 \pm 0.030 \\ 0.259 \pm 0.077 \\ 0.214 \pm 0.034 \\ 0.223 \pm 0.027 \end{array}$	$\begin{array}{c} 0.523 \pm 0.027 \\ 0.287 \pm 0.114 \\ 0.212 \pm 0.037 \\ 0.217 \pm 0.028 \end{array}$	$\begin{array}{c} \underline{0.549} \pm 0.006 \\ \underline{0.303} \pm 0.085 \\ \hline \textbf{0.240} \pm 0.026 \\ 0.244 \pm 0.025 \end{array}$	$\begin{array}{c} 0.539 \pm 0.026 \\ 0.299 \pm 0.089 \\ \underline{0.240} \pm 0.030 \\ \underline{0.247} \pm 0.025 \end{array}$	0.546 0.273 0.218 0.230
spam_base	С	Accuracy ↑ Precision ↑ F1 ↑ Recall ↑	$\begin{array}{c} \textbf{0.939} \pm 0.008 \\ \textbf{0.939} \pm 0.007 \\ \underline{0.929} \pm 0.010 \\ \textbf{0.928} \pm 0.008 \end{array}$	$\begin{array}{c} 0.927 \pm 0.017 \\ 0.930 \pm 0.017 \\ 0.921 \pm 0.018 \\ 0.915 \pm 0.019 \end{array}$	$\begin{array}{c} 0.929 \pm 0.004 \\ 0.932 \pm 0.006 \\ 0.925 \pm 0.005 \\ 0.919 \pm 0.005 \end{array}$	$\begin{array}{c} 0.933 \pm 0.009 \\ 0.935 \pm 0.009 \\ 0.927 \pm 0.009 \\ 0.921 \pm 0.009 \end{array}$	$\frac{0.936}{0.939} \pm 0.007$ $\frac{0.939}{0.932} \pm 0.008$ 0.932 ± 0.007 0.927 ± 0.008	0.912 0.919 0.906 0.898

Table 5: Overall performance comparison. The best results are highlighted in **bold**, and the second-best results are <u>underlined</u>. (↑ indicates that a higher value of the metric corresponds to better performance, while ↓ denotes that a lower value of the metric indicates better performance.)

834

835

836

837

evaluator with LR, DT, GBDT, RF and XGB respectively. Table 5 and Table 2 shows the comparison results in terms of overall performance. EASE as an evaluator, significantly enhances feature selection performance, demonstrating outstanding capability. This indicates that EASE not only possesses excellent generalization ability but also excels in evaluation performance. The potential reason for this superior performance lies in the Feature-Sample Subspace Generator, which greatly improves generalization, while the Contextual Attention Evaluator further optimizes performance by capturing interactions within the feature space.

842

C.3 ADDITIONAL EFFICIENCY COMPARISON

We additionally compare training time across all datasets. We choose RFE for iterative feature 843 selection and adopt EASE, GBDT, RF and XGB as the feature space evaluator respectively. Figure 6 844 shows the comparison results in terms of cumulative time. And we have omitted algorithms that 845 consume too much time. We find that EASE outperforms other baselines nearly across all datasets. 846 Specifically, EASE can save over 100 seconds of evaluation time compared to other time-consuming 847 baselines on bike_share, wine_white, openml_616, openml_607, and openml_586. The underlying 848 driver is that our incremental parameter update strategy focuses on the most relevant information and 849 quickly capture evolving patterns in the feature space for evaluation, leading to a faster optimization 850 process. In conclusion, EASE can effectively and efficiently evaluate the quality of the feature 851 space.

852 853

854

C.4 THE EFFECTIVENESS OF EASE FOR FEATURE SPACE REFINEMENT

855 We additionally test the prediction performance on all datasets between the original feature space 856 and the space refined by EASE using various downstream predictors, including LR, DT, GBDT, and 857 RF. Figure 7 shows the overall comparison results in terms of Accuracy and MAE according to the 858 task type. We find that the feature space evaluated by EASE outperforms the original feature space 859 across all datasets and baselines. For the datasets openml_616, openml_607, spam_base, spectf, and 860 svmguide3, and wine_white, EASE significantly improves performance in terms of Accuracy or 861 MAE. The underlying driver is that our information decoupling strategy can effectively integrate information interactions and provide it to the contextual attention evaluator. Then contextual attention 862 evaluator accurately captures the intrinsic interactions of feature space, thereby guiding the feature 863 space iterative optimization algorithm to obtain high-quality feature space.

EASE 0.8 0.0 (e) airfoil (a) openml_616 (b) openml_607 (c) openml_586 (d) openml_620 EASE EASE (g) spectf (f) spam_base (h) symguide3 (i) wine red (j) wine_white



C.5 **EASE'S PERFORMANCE IN DIFFERENT FEATURE SELECTION FRAMEWORKS**

884 This experiment aims to answer: Is EASE generalizable and applicable across different feature 885 space optimization algorithms? We apply EASE to a state-of-the-art FS algorithm SDAE (Has-886 sanieh & Chehade, 2024). SDAE learns low-dimensional representations of high-dimensional data 887 through a deep auto-encoder structure, while introducing a selective layer that automatically selects 888 a relevant subset of features representing the entire feature space. This method performs FS in an 889 unsupervised manner, effectively capturing nonlinear relationships between features. We respectively selected EASE LR, DT, GBDT, RF, and XGB as evaluators to assess the the final selected 890 feature space. Table 6 shows the comparison results across different datasets. We can observe that in 891 the feature space selected by the SDAE algorithm, EASE exhibits the best evaluation performance 892 in both classification and regression tasks. This proves that EASE has stronger robustness and 893 can accurately capture the key information in the feature space compared with other methods. The 894 underlying reason for this lie in EASE's innovative design of the Contextual Attention mechanism. 895

896 Table 6: Comparison of different evaluators in terms of Accuracy (for classification tasks) and MAE 897 (for regression tasks) in SDAE framework. The best results are highlighted in **bold**. The second-best 898 results are highlighted in underline.(Lower MAE value and higher Accuracy value corresponds to 899 better performance.

~								
0	Dataset	R/C	EASE	LR	DT	GBDT	RF	XGB
	openml 607	R	0.232 ± 0.034	0.740 ± 0.040	0.397 ± 0.030	0.233 ± 0.015	0.278 ± 0.012	0.243 ± 0.011
2	openml_616	R	0.202 ± 0.014	0.764 ± 0.045	0.518 ± 0.024	$\frac{0.233}{0.314} \pm 0.026$	0.343 ± 0.030	0.331 ± 0.033
	openml_620	R	$\textbf{0.199} \pm 0.005$	0.735 ± 0.033	0.477 ± 0.007	$\overline{0.327} \pm 0.013$	0.358 ± 0.024	0.324 ± 0.012
	openml_586	R	$\underline{0.228} \pm 0.016$	0.713 ± 0.041	0.383 ± 0.008	$\textbf{0.227} \pm 0.019$	0.264 ± 0.019	$\overline{0.229} \pm 0.010$
	airfoil	R	$\overline{\textbf{0.154}} \pm 0.013$	0.564 ± 0.025	0.273 ± 0.012	0.286 ± 0.017	0.195 ± 0.008	$\underline{0.162} \pm 0.012$
	bike_share	R	$\textbf{0.007} \pm 0.000$	0.020 ± 0.000	0.016 ± 0.001	0.019 ± 0.000	0.007 ± 0.001	0.013 ± 0.001
	wine_red	С	$\textbf{0.741} \pm 0.037$	0.591 ± 0.008	0.549 ± 0.022	0.640 ± 0.020	0.724 ± 0.032	0.600 ± 0.009
	svmguide3	С	0.856 ± 0.014	0.806 ± 0.021	0.795 ± 0.018	0.832 ± 0.024	$\textbf{0.861} \pm 0.012$	0.851 ± 0.004
	wine_white	С	$\textbf{0.686} \pm 0.034$	0.518 ± 0.011	0.585 ± 0.010	0.584 ± 0.009	0.675 ± 0.015	0.654 ± 0.008
	spectf	С	0.829 ± 0.046	0.789 ± 0.064	0.772 ± 0.030	0.756 ± 0.072	$\textbf{0.837} \pm 0.057$	0.813 ± 0.023
	mammography	С	$\textbf{0.992} \pm 0.001$	0.983 ± 0.001	0.980 ± 0.003	0.986 ± 0.002	0.987 ± 0.004	0.987 ± 0.001
	spam_base	С	$\textbf{0.963} \pm 0.006$	0.927 ± 0.007	0.904 ± 0.013	0.946 ± 0.003	0.960 ± 0.007	0.946 ± 0.005
	AmazonEA	С	$\textbf{0.950} \pm 0.003$	0.944 ± 0.002	0.930 ± 0.003	0.942 ± 0.006	$\underline{0.947} \pm 0.002$	0.943 ± 0.001
	Nomao	С	$\textbf{0.973} \pm 0.002$	0.941 ± 0.004	0.945 ± 0.003	0.953 ± 0.002	0.967 ± 0.001	0.969 ± 0.002

912 913

914

864

865 866

871

872

879 880

881 882

883

C.6 THE IMPACT OF EACH TECHNICAL COMPONENT

915 This experiment aims to answer: How does each technical component in EASE impact its efficiency? We compare the average training time of EASE with other EASE variants in each opti-916 mization, including EASE^{-FC}, EASE^{-IT} and EASE^{-PT}. We develop EASE^{-PT}, EASE^{-IT}, 917 and $EASE^{-FC}$ by removing the pre-training, incremental training, and feature-sample subspace 918 construction steps from EASE respectively. Table 7 shows the comparison results across different 919 datasets in terms of average training time and standard deviation. We observe that EASE has the 920 shortest runtime across all the datasets. Specifically, for classification, EASE-IT has the second-921 best performance which indicate the technical component of feature subspace construction and pretraining are crucial for enhancing efficiency. For regression, $EASE^{-PT}$ exhibit the second highest 922 efficiency which demonstrate the technical component of incremental training and feature subspace 923 construction can significantly improve EASE evaluation efficiency. The possible reason is pre-924 training can obtain a well-initialized contextual attention evaluator to provide a strong foundation 925 for evaluation. Incremental training can leverage the overlapping information between consecutive 926 iterations to avoid redundant computations and training time. And feature subspace construction 927 can decouple the information within the feature space. In summary, our proposed EASE, which 928 includes components for pre-training, incremental training, and feature subspace construction sig-929 nificantly reduce the average training time. 930

Table 7: Comparison of different EASE variants in terms of time complexity. The best results are highlighted in **bold**. The second-best results are highlighted in <u>underline</u>. (The unit is seconds.)

Dataset	R/C	EASE	$EASE^{-FC}$	$EASE^{-IT}$	$EASE^{-PT}$
openml_607	R	$\textbf{4.396} \pm 1.237$	5.400 ± 0.134	5.762 ± 0.044	5.175 ± 0.076
openml_616	R	$\textbf{3.994} \pm 1.347$	5.408 ± 0.639	5.859 ± 0.051	$\overline{5.097} \pm 0.074$
openml_620	R	$\textbf{4.547} \pm 1.322$	$\underline{5.635}\pm0.063$	6.598 ± 0.068	5.687 ± 0.142
openml_586	R	$\textbf{4.644} \pm 0.382$	7.287 ± 0.271	7.459 ± 0.056	6.980 ± 0.212
airfoil	R	$\textbf{4.201} \pm 0.595$	6.206 ± 0.340	6.531 ± 0.064	5.792 ± 0.340
bike_share	R	$\textbf{3.356} \pm 0.069$	5.877 ± 0.122	6.501 ± 0.116	5.745 ± 0.116
wine_red	С	$\textbf{2.848} \pm 0.273$	5.998 ± 1.006	3.235 ± 0.123	4.585 ± 0.062
svmguide3	С	$\textbf{0.872} \pm 0.134$	5.899 ± 0.272	1.538 ± 0.468	6.178 ± 0.540
wine_white	С	$\textbf{3.815} \pm 1.228$	6.180 ± 0.225	$\underline{4.083} \pm 0.671$	6.002 ± 0.103
spam_base	С	$\textbf{4.032} \pm 0.354$	7.598 ± 1.035	4.745 ± 0.829	7.622 ± 1.046
mammography	С	$\textbf{0.733} \pm 0.010$	1.436 ± 0.024	0.813 ± 0.088	4.799 ± 0.020
spectf	С	$\textbf{3.005} \pm 0.043$	5.852 ± 0.577	$\underline{3.216}\pm0.301$	5.276 ± 0.343

944 945 946

947

C.7 PARAMETER SENSITIVITY ANALYSIS

948 This experiment aims to answer: How do parameters affect the performance of EASE? To validate 949 the parameter sensitivity of key parameters in EASE, we select the wine_white and openml_586 950 datasets as examples. We focus on the number of heads h and the embedding dimension D in 951 training procedure. To address the issue of varying feature space lengths during the iterative pro-952 cess, we first set the size of the feature subspace to match the embedding dimension and then 953 transpose it, successfully overcoming this challenge. Consequently, D is both the embedding di-954 mension and the size of the feature subspace. Specifically, We set D = 32 and test the value 955 of h with the set $\{2, 4, 8, 16, 32\}$. And we set h = 16 and test the value of D with the set $\{16, 32, 48, 64, 80, 96, 112, 128\}$. Figure 8 shows the comparison results in terms of Accuracy, Re-956 call and F1 Score for classification task, 1-MAE, 1-RMSE, and R^2 Score for regression task. 1-MAE 957 and 1-RMSE used for denoting that a higher value of the metric indicates better performance. We 958 observe that the performance of downstream tasks generally remains stable across different values 959 of h and D, with significant changes occurring only at specific parameter values, such as h = 4 and 960 D = 64 for the regression. A possible reason for this observation is that our proposed EASE can 961 effectively decouple information within feature space and can capture contextual information during 962 evaluation process. This observation indicates that EASE is not sensitive to the number of heads h963 and the embedding dimension D. Therefore, the evaluating process of EASE is robust and stable.

964 965 966

C.8 CASE STUDY

967 This experiment aims to answer: What is the difference between the original feature space and the 968 feature space refined by the EASE for ML tasks? We select the wine_white dataset as example to 969 visualize its features. In detail, we use the original feature space and a refined feature space evaluated 970 by EASE within the RFE framework, with RF as the downstream predictor. Figure 9 shows the 971 importance of top 8 features and their impact on the original feature space and EASE feature space. 972 Specifically, we select 400 samples and calculate their contribution during the prediction process.



(a) D on openml_586 (R). (b) h on openml_586 (R). (c) D on wine_white (C). (d) h on wine_white (C).

Figure 8: Parameter sensitivity on the number of heads h and the embedding dimension D on wine_white and openml_586 datasets.

Table 8: Comparison of different evaluators in terms of Accuracy (for classification tasks) and MAE (for regression tasks) in GRFG framework. The best results are highlighted in **bold**. The second-best results are highlighted in <u>underline</u>.

Dataset	R/C	EASE	LR	DT	GBDT	RF	XGB
openml_616	R	$\textbf{0.302} \pm 0.006$	0.321 ± 0.002	0.327 ± 0.006	0.331 ± 0.012	$\underline{0.312}\pm0.012$	0.338 ± 0.000
openml_586	R	$\textbf{0.373} \pm 0.001$	0.410 ± 0.008	0.410 ± 0.009	0.409 ± 0.003	0.385 ± 0.017	0.399 ± 0.002
svmguide3	С	$\textbf{0.849} \pm 0.002$	0.816 ± 0.000	0.821 ± 0.006	0.822 ± 0.006	$\overline{0.826} \pm 0.000$	0.820 ± 0.005
mammography	С	$\textbf{0.993} \pm 0.001$	0.984 ± 0.000	0.985 ± 0.018	0.986 ± 0.001	$\overline{0.985}\pm0.000$	$\underline{0.986} \pm 0.000$

981

982

983 984

985

986

The horizontal axis represents the SHAP values for each feature, reflecting the impact of that feature 994 on the prediction, while the vertical axis lists the feature names in order of importance (Temenos 995 et al., 2023). And we color the size of the feature values (red represents larger values, while blue 996 represents smaller values). We find that the EASE feature space greatly enhances the predictor 997 accuracy by 15%. Another interesting observation is that the feature ranking in the EASE feature 998 space differs from that in the original feature space for the same predictor RF. In detail, we can trace 999 and explain the source and effect of specific feature. For example, "volatile acidity" measures the 1000 impact of the wine's acidity on the wine quality evaluation, which is positively correlated with wine 1001 quality. The underlying driver for these observations is that the multi-head attention mechanism 1002 in contextual attention evaluator can capture the interactions between samples and features after 1003 decoupling the information, which not only improves the fairness of the evaluation but also enhances its interpretability. Thus, this case study reflects that the effectiveness and interpretability of EASE 1004 as a evaluator for feature space quality evaluation. 1005

For all other details of the hyperparameter configurations, optimization strategies, specific training processes, and environmental settings, please refer to Appendix C.1.



1018 1019

1020

Figure 9: Comparison of feature importance in EASE feature space and original feature space.

(b) Original Feature Space (Accuracy = 0.4215).

1021 1022 1023

1024

C.9 EASE's Performance in Feature Generation Frameworks

(a) EASE Feature Space (Accuracy = 0.5730).

1025 This experiment aims to answer: *Is* EASE *generalizable and applicable in feature generation algorithm?* We apply EASE to a feature generation method GRFG (Wang et al., 2022). GRFG addresses challenges in representation space reconstruction by proposing a cascading deep rein-forcement learning approach that automates feature generation through a group-wise strategy and nested interactive processes. We respectively selected EASE LR, DT, GBDT, RF, and XGB as eval-uators to evaluate the performance of the generated or selected features during the GRFG procedure. Table 8 shows the comparison results across different datasets. We observe that the proposed EASE achieves the best performance in both classification and regression tasks. This further demonstrates that, compared to traditional feature evaluation algorithms, EASE exhibits excellent performance in both feature selection and feature generation tasks and can capture the key information of the feature space.