DISTRIBUTION SHIFT AWARE NEURAL FEATURE TRANSFORMATION

Anonymous authors

Paper under double-blind review

Abstract

Feature transformation, as a core task of Data-centric AI (DCAI), aims to improve the original feature set to enhance AI capabilities. In dynamic real-world environments, where there exists a distribution shift, feature knowledge may not be transferable between data. This matter prompts a distribution shift feature transformation (DSFT) problem. Prior research works for feature transformation either depend on domain expertise, rely on a linear assumption, prove inefficient for large feature spaces, or demonstrate vulnerability to imperfect data. Furthermore, existing techniques for addressing the distribution shift cannot be directly applied to discrete search problems. DSFT presents two primary challenges: 1) How can we reformulate and solve feature transformation as a learning problem? and 2) What mechanisms can integrate shift awareness into such a learning paradigm? To tackle these challenges, we leverage a unique Shift-aware Representation-Generation Perspective. To formulate a learning scheme, we construct a representation-generation framework: 1) representation step: encoding transformed feature sets into embedding vectors; 2) generation step: pinpointing the best embedding and decoding as a transformed feature set. To mitigate the issue of distribution shift, we propose three mechanisms: 1) shift-resistant representation, where embedding dimension decorrelation and sample reweighing are integrated to extract the true representation that contains invariant information under distribution shift; 2) flatness-aware generation, where several suboptimal embeddings along the optimization trajectory are averaged to obtain a robust optimal embedding, proving effective for diverse distribution; and 3) shift-aligned pre and post-processing, where normalizing and denormalizing align and recover distribution gaps between training and testing data. Ultimately, extensive experiments are conducted to indicate the effectiveness, robustness, and trackability of our proposed framework. Our code is available at https://tinyurl.com/OODFT.

034 035

037

004

006

008 009

010 011

012

013

014

015

016

017

018

019

021

024

025

026

027

028

029

031

032

1 INTRODUCTION

In the modern era of machine learning, models
are large and effective yet GPUs are expensive.
Instead, Data-centric AI (DCAI) has emerged
as an alternative solution by using AI to augment data power for better AI, even with simple models (Zha et al., 2023). As a core task of
DCAI, feature transformation aims to transform



Figure 1: An illustration of DSFT, where the distribution of the training and test sets differ.

an original feature set into a more powerful one (Kusiak, 2001; Ying et al., 2023). In most cases,
we assume that training and testing data are from the same domains and are I.I.D. without distribution shifts. In open environments, feature distributions, even from the same domain, can shift from training to testing data. Feature knowledge learned from data before shifts may not apply to data after shifts (see Figure 1). This case can be generalized as a distribution shift feature transformation (DSFT) problem. Addressing DSFT can enhance the generalization, robustness, applicability, and availability of feature transformation.

Relevant works can only partially solve DSFT. Firstly, DSFT is related to feature transformation.
 For instance, humans can manually transform a feature set with domain knowledge and empirical experiences. Machine-assisted methods include principal components analysis (PCA) (Maćkiewicz

054 & Ratajczak, 1993), exhaustive-expansion-reduction approaches (Kanter & Veeramachaneni, 2015; 055 Khurana et al., 2016), iterative-feedback-improvement approaches (Khurana et al., 2018; Tran et al., 056 2016). However, human-based methods are usually time-consuming and incomplete. Among 057 machine-assisted methods, PCA is only based on a straight linear feature correlation assumption, 058 ignoring non-linear feature interaction, and PCA only reduces dimensionality instead of increasing dimensionality; others are mostly based on a discrete search formulation with critical limitations: 1) 059 large search space and time costly when finding optimal, 2) lacks computational generalization and 060 robustness against imperfect data, such as distribution shifts. Second, DSFT is related to anti-shift 061 learning, generalization, and adaptive learning, for instance, distribution alignment (Fan et al., 2024; 062 Kim et al., 2021; Hu et al., 2023), parameter isolation (Zhang et al., 2023), adaptive experience 063 replay (Li et al., 2024). A clear research gap is that it is difficult to integrate shift awareness in 064 learning-based formulation into discrete search-based formulations. 065

There are two major challenges to bridge the gap between feature transformation and shift aware-066 ness: 1) From search-based formulation to learning-based formulation. Many existing studies 067 view feature transformation as a task of searching for the best feature set among large discrete fea-068 ture combination possibilities. The first challenge seeks to answer: what paradigm can solve feature 069 transformation as a learning problem: within a continuous optimization space, measured by a discriminative or generative model, evaluated by a tangible objective, guided by directional gradient? 071 2) From shift-sensitive to shift-robust. Once the feature transformation learning paradigm is pro-072 posed, the second challenge intends to answer: what mechanisms can integrate shift awareness into 073 such a learning paradigm?

074 Our Insights: A Shift-aware Representation-Generation Perspective. We formulate the DSFT 075 problem into a deep learning task, instead of a discrete search task. The study in (Wang et al., 076 2023) found that by collecting transformed feature sets and task performances as training data, 077 we can embed transformed feature sets as embedding vectors. This continuous embedding space 078 of transformed feature sets, if accurate, should include and describe the embedding point of the 079 unobserved best transformed feature set at a certain maximum. Guided by this finding, we propose to regard searching the best feature transformation as a gradient optimization problem in the embedding 081 space, in order to identify the optimal embedding. To address the learning framework challenge, we develop a representation-generation framework: 1) the representation step encodes transformed feature sets into embedding vectors; 2) the generation step identifies the best embedding via gradient 083 ascent and decodes the best embedding to the best transformed feature set. Most importantly, this 084 framework, although widely used in computer vision, opens up three essential opportunities for 085 integrating shift awareness into representation, generation, and pre/post-processing for an old-school 086 feature engineering problem. 087

880 To address the distribution shift, we develop the following three mechanisms in response to shiftresistant representation, flatness-aware generation, and shift-aligned pre and post-processing. (i) 089 Connecting embedding dimension decorrelation with sample reweighing can incorporate shift re-090 sistance into embedding. (ii) Leveraging flatness in gradient-based optimal embedding search can 091 ensure that the performance of a transformed feature set embedding does not decrease significantly 092 in a neighborhood around the maximum, thus mitigating distribution shift. (iii) Performing normal-093 ization in pre-processing and denormalization in post-processing can align distribution gaps between 094 training and testing data. 095

Summary of Proposed Solution. Inspired by these insights, we propose a principled 096 representation-generation-based deep feature transformation learning framework, where the representation step trains an encoder to encode transformed feature sets into embedding vectors, by 098 optimizing an evaluator (whether the embedding of a transformed feature set can be used to predict its corresponding performance) and a decoder (whether a transformed feature set can be recon-100 structed given its corresponding embedding), and the generation step exploits gradient ascent to 101 identify the best transformed feature set embedding, and use the well-trained decoder to decode the 102 best embedding and generate the best transformed feature set. This framework's significance lies 103 not in its novelty, but in its ability to offer flexibility, thereby facilitating awareness shifts across three channels. Specifically, to achieve shift-resistant embedding, we reformulate a feature set as 104 a feature-feature similarity graph and develop a new graph neural network that jointly considers 105 partial covariance matrix norm minimization, sample reweighing in evaluator loss residuals, and 106 decoder optimization, and is optimized by a bilevel training strategy. To achieve flatness-aware fea-107 ture transformation generation, we connect oscillating search and suboptimal embedding averaging

to achieve a flat region in optimal embedding identification for the robust transformed feature set
 decoding. Finally, we leverage normalization and denormalization in pre and post-processing to
 further alleviate shifts.

Our contributions are: 1) we formulate DSFT as a continuous optimization problem and employ a representation-generation scheme to identify the optimal transformed feature space; 2) we propose shift-resistant feature set representation, flatness-aware generation, and integration of normalization to address distribution shift; 3) we conduct extensive experiments to demonstrate the efficiency, resilience, and traceability of our framework.

116 117 118

2 PROBLEM STATEMENT

119 The DSFT Problem. Given a feature set, a training dataset, and a test dataset, and both training and 120 test data share the same feature set, considering the existence of distribution shifts between training and test data samples, we aim to learn an anti-shift feature transformation model to transform the 121 original shifted training and test datasets into a transformed training dataset and a transformed test 122 dataset, in order to: 1) improves the performance of an ML task (e.g., regression, classification); 2) 123 is anti-shift. Our perspective is that the knowledge about what feature structure advances data power 124 can be learned from historical exploration experiences of feature transformations and performances. 125 We convert the DSFT problem as a shift-aware generative feature transformation learning task that 126 involves several concepts: 127

1) Feature Cross. We apply operations to features, generating feature cross (e.g., $f_1 + f_1/f_3 - f_2$).

2) Feature Transformation Operation Sequence. We combine multiple feature crosses to construct a new transformed feature set (e.g., $f_1 + f_1/f_3 - f_2, \sqrt{f_1}$).

3) The postfix expression of feature transformation operation sequence. We use a postfix expression to represent the feature transformation operation sequence. For instance, the postfix expression of f_1+f_1/f_3-f_2 , $\sqrt{f_1}$ is: $SOSf_1f_1f_3/+f_2-SEPf_2\sqrt{.EOS}$. The postfix expression reduces token redundancy, prevents illegal transformation operations and semantic ambiguity, and minimizes search space. After converting a transformed feature set into a postfix token expression of a feature transformation operation sequence, DSFT task can be solved as a sequential generation task that generates a feature transformation operation sequence.

138 139

140

3 RL AGENTS FOR DATA COLLECTION

Before modeling, we collect various transformed feature sets and their corresponding performances (e.g., accuracy) on a specific downstream task (e.g., regression, classification) as training data.

143 Why automated and diverse training data collection. To learn a feature set2vec encoder, we need 144 to collect various feature sets and their accuracy on a downstream task as a supervised knowledge 145 base. Intuitively, we can manually and randomly try different possible combinations of feature crosses (e.g., $f_1 + f_2, f_1/f_2$. Here, f_1 is a head feature selected by a head agent, and f_2 is a tail 146 feature selected by a tail agent. f_1 and f_2 can represent any feature.) to form a new feature set, 147 and then test the performance of the new feature set on a controlled regression or classification task. 148 However, the three aspects of training data are essential: 1) large volume: sufficient training data; 2) 149 high quality: high-accuracy feature set cases as successful experiences; 3) high diversity: training 150 data should not ignore random, exploratory, and failure cases. 151

Leveraging reinforcement learning to ex-152 plore high-quality and diverse training data. 153 Inspired by the exploitation, exploration, and 154 self-learning abilities of reinforcement learn-155 ing, our idea is to view an RL agent as a training 156 data collector in order to achieve volume (self-157 learning enabled automation), diversity (explo-158 ration), and quality (exploitation). Specifically, 159 we design RL agents to automatically decide



Figure 2: Data preparation pipeline. RL agents are employed to collect exemplary and diversified training data.

how to perform feature cross and generate new feature sets. The reinforcement exploration experi how to perform feature cross and generate new feature sets. The reinforcement exploration experi and corresponding accuracy will be collected and stored as training data as shown in Figure 2.
 More details about the RL data collector are described in Appendix A.

162 4 THE REPRESENTATION-GENERATION FRAMEWORK

164 4.1 FRAMEWORK OVERVIEW

165 Figure 3 shows our Shift Awareness Feature 166 Transformation framework (SAFT) includes 167 three components: 1) shift-resistant feature 168 set representation; 2) flatness-aware feature transformation generation; 3) integrated pre-170 normalization and post-denormalization. То achieve shift-resistant representation, we de-171 velop two insights: 1) feature sets as dynamic 172 feature-feature interaction graphs and 2) a com-173 bination of sample reweighing and embedding 174 dimension orthogonality for invariant represen-175 tation. To achieve flatness-aware feature trans-176 formation generation, we integrate flatness-177 aware gradient ascent, which aims to identify a 178 flatter optimal embedding for shift robustness. 179 Finally, we integrate normalization in prepro-



Figure 3: An overview of our framework.

cessing and denormalization in postprocessing to further alleviate the impacts of shifts. The Algorithm 1 in Appendix B.1 provides the detailed pseudo-code.

182 183

4.2 Shift-resistant Feature Graph Embedding

To find the best transformed feature sets, our key insight is to effectively represent feature sets 184 and corresponding performances in an embedding space. Using training data composed of feature 185 sets and their predictive performances, we train an encoder and a decoder. These components map transformed feature sets to vectors and vice versa, creating a continuous embedding space. We 187 believe the best transformed feature set, termed feature knowledge, is described in this continuous 188 embedding space and, thus, can be efficiently identified by gradient optimization. There is no need 189 to explore exponentially growing possibilities of feature combinations. Although training-testing 190 shifts may distort the embedding space, the space is learnable. Anti-shift mechanisms can rectify 191 and adjust the distortion and ensure reliable and transferable feature knowledge in the testing set. 192

A. Encoder for Learning Dynamic Feature-Feature Interaction Graph Embedding. Representing a discrete feature set is traditionally challenging. Conventional methods either aggregate the first-order or second-order statistics of each feature or regard a set of features as a token sequence for encoding. However, a feature set is not a decorrelated independent set. Features within a set strongly associate to demonstrate distinctive patterns and enhance ML performances. When ignoring feature-feature interaction, feature set representations are likely to be biased and inaccurate, which introduces errors in the optimal feature set identification.

Our insight is to view a feature set as a feature-feature interaction attributed graph, where a feature 200 is a node, the element values of a feature are node attributes, and feature-feature similarity (i.e., 201 cosine similarity) is the weight of an edge. After converting feature sets into graphs, existing graph 202 embedding methods, for example, Graph Convolutional Networks (GCNs), could have been applied 203 to map a feature set to an embedding vector. However, we observe that the sizes (feature numbers) 204 and topology (adjacency matrices) of feature-feature similarity graphs vary a lot. Classic GCNs 205 are unable to address dynamic graph structure (sizes and topology) because GCNs require fixed 206 graph sizes and topology. Inspired by (Hamilton et al., 2017), we tailor the idea of node sampling and neighbor information propagation from increasingly distant reaches for dynamic feature-feature 207 graph embedding. Specifically, given a feature-feature similarity attributed graph at each step, each 208 node consolidates the previous representations of itself and its neighbor representations. A fully con-209 nected layer with a nonlinear activation function is later utilized to update the current representation. 210 We normalize the embedding for each step and iterate to obtain the final representation. 211

B: Estimator for Estimating Feature Set Performances. Our goal is to learn an embedding space to
represent feature sets so that we can find the optimal feature set embedding with higher performance,
which is used to reconstruct the optimal feature set. To strengthen the performance expressiveness of
feature set embedding, we incorporate a feature set performance estimator as a downstream task of
the feature-feature graph encoder. Specifically, we employ a feedforward network to take a feature-

feature graph embedding as input and predict the performance of a feature set, by minimizing the mean square error between estimated feature set performances and ground-truth feature set performances. Formally, the loss of evaluator is $\mathcal{L}_{est} = \frac{1}{N} \sum_{i=1}^{N} (p_i - \omega_{\theta}(E_i))^2$, where N is total number of feature sets in training data, *i* index each feature set, p_i is the ground-truth feature set performance, ω_{θ} is the feature set performance estimator given a feature set's graph embedding, E_i is the sample embedding.

222 C: Decoder. The primary objective of the decoder is to output a feature transformation operation 223 sequence (deemed as a token sequence) given a feature-feature graph embedding. The decoder 224 provides two main functionalities: 1) enforcing the faithfulness of embedding: the embedding of 225 a transformed feature set should be able to fully reconstruct its corresponding feature transforma-226 tion operation sequences; 2) generates the optimal transformed feature set: a well-trained decoder 227 can be used to generate the feature transformation operation sequence of the optimal embedding. 228 Our decoder is structured by a single-layer LSTM and a softmax layer, where the LSTM is to recursively learn the hidden state over the current-previous dependent feature transformation op-229 eration sequence, and the softmax layer is to estimate the token probability when generating each 230 token. Specifically, given ψ as a single-layer LSTM, the distribution of the single *i*-th token is 231 $P_{\psi}(\gamma_i|E, \Upsilon_{<i}) = \frac{\exp(s_j)}{\sum_M \exp(s)}$. where, Υ is a feature transformation operation token sequence with M token length, γ_i is the *i*-th token in Υ , and s_j is the *j*-th output of the softmax layer. Finally, we aim to minimize the negative log-likelihood of decoding a feature transformation operation token 232 233 234 sequence as the reconstruction loss: $\mathcal{L}_{rec} = -\sum_{N} \sum_{i=1}^{M} log P_{\psi}(\gamma_i | E, \Upsilon_{< i}).$ 235

D: Shift-resistant Joint Optimization of Encoder, Estimator, and Decoder. We will develop a shift-resistant optimization strategy to learn the encoder, estimator, and decoder.

The Joint Objective Function. The separate training of the encoder, estimator, and decoder can lead to local optima. Instead, we optimize the joint loss of the encoder, evaluator, and decoder. Our objective is: $\mathcal{L}_{est} + \gamma \mathcal{L}_{rec}$, where γ are a hyper-parameter. In this way, we enforce the encoderevaluator-decoder structure to learn a faithful embedding to accurately predict the performance of a transformed feature set and decode the corresponding feature transformation operation sequence.

Debiasing Shift-induced Biases in Training for Shift-aware Learning. In open environments, 244 there exist distribution shifts between training and testing data. The shifts introduce bias into the 245 embedding space learned from training data, which can not be generalized to shifted test data. The 246 learned biased embedding thus includes not just true representations, but also false representations. 247 The true representations are relevant and genuine descriptors of the performance of a transformed 248 feature set, and remain invariant under distribution shifts. The false representations are irrelevant 249 descriptors that drift from training to testing. The subtle and fake correlations between false rep-250 resentations and true representations mislead the joint loss optimization process to relate false rep-251 resentations to the labels of feature set performances, thus, degrading generalization abilities under 252 distribution shifts.

253 To alleviate the negative influence of false representations, an intuitive idea is to decorrelate true and 254 false representations during learning. However, true and false representations are unknown and not 255 pre-identified. We are unable to directly decorrelate true and false representations. A conservative 256 and aggressive solution is to decorrelate all representations by encouraging the orthogonality of all 257 embedding dimensions. This objective can be reformulated into the minimization of the squared 258 Frobenius norm of a partial covariance matrix of the embedding dimensions. This Frobenius norm 259 can be added as a loss term into the joint loss minimization. However, 1) it relies on the aggressive 260 and strong assumption on the structure of embedding dimensions, thus, suffers from loss of information; 2) even though all embedding dimensions are decorrelated, the false representations are 261 still in the embedding. In other words, feature set embeddings are still not pure and suffer from the 262 inclusion of false and biased dimensions, resulting in low generalization over changing distributions. 263

Our goal is to debias the embedding bias caused by distribution shift so the embeddings include the true invariant correlations between true representations and feature set performances. We connect and unify two interesting insights: 1) A recent study (Shen et al., 2020) shows that there exists a set of sample weights that can reshape the sample-variable matrix and make a weighted covariate distribution matrix near orthogonal. 2) Another study (Athey et al., 2018) shows that reweighing the losses of training samples of transformed feature sets can achieve a debiasing effect in selecting variables to rebuild representations for better generalization. Interestingly, this insight closely 270 connects to the boosting strategy that reweighs training samples (i.e., increase misclassified sample 271 weights and lower correctly classified sample weights) to enforce the next weak classifier to avoid 272 making mistakes in previous misclassified samples. In summary, training sample reweighing can 273 be seen as model regularization. Inspired by the two findings and (Zhang et al., 2021; Li et al., 274 2022), we develop a unified bilevel optimization approach: 1) the inner-loop level: we see the sample weights as learnable parameters that are learned by minimizing the squared Frobenius norm of 275 a partial covariance matrix; 2) the outer-loop level: we use the weights to reweigh the sample losses 276 in estimating the performance of each feature-feature graph embedding. In this way, we leverage 277 sample reweighing and bi-level training as a bridge to achieve both decorrelating all embedding 278 dimensions and reducing false dimensions caused by the bias of distribution shift. 279

The Shift-resistant Bilevel Training. The Inner-loop Training. For the inner-loop, our objective is to learn the best sample weights to minimize the squared Frobenius norm as follows.

280

281

284

 $\mathbf{R}^* = \underset{\mathbf{R}}{\operatorname{arg\,min}} \sum_{i < j} \|\hat{\mathbf{C}}_{E_{*i}, E_{*j}}^{\mathbf{R}}\|_F^2, \tag{1}$

where $\hat{\mathbf{C}}_{E_{*i},E_{*j}}^{\mathbf{R}} = \frac{1}{N-1} \sum_{n=1}^{N} [(r_n f(E_{ni}) - \frac{1}{N} \sum_{m=1}^{N} r_m f(E_{mi}))^\top \cdot (r_n g(E_{nj}) - \frac{1}{N} \sum_{m=1}^{N} r_m g(E_{mj}))]$. Here, $\hat{\mathbf{C}}_{E_{*i},E_{*j}}^{\mathbf{R}}$ is the partial cross-covariance matrix, $\mathbf{R} = \{r_n\}_{n=1}^{N}$ is the graph weight vector, r_i is the weight of *i*-th feature subset graph and we constraint $\sum_{n=1}^{N} r_n = N$. $f(\cdot)$ and $g(\cdot)$ are the random Fourier features functions, E_{*i} and E_{*j} are the different dimensions of the same training sample.

The Outer-loop Training. We use the updated weights of training samples $\mathbf{R}^* = \{r_n^*\}_{n=1}^N$ to reweigh the estimation loss of each transformed feature set, given by $\mathcal{L}_{est} = \sum_{i=1}^N r_i^* (p_i - \omega_{\theta}(E_i))^2$. Then, we combine the new estimator loss and decoder loss \mathcal{L}_{rec} to obtain the weighted joint loss: $argmin \ \mathcal{L} = \mathcal{L}_{est} + \gamma \mathcal{L}_{rec}$. The Algorithm 2 in Appendix B.2 provides the detailed pseudo-code.

296 4.3 FLATNESS-AWARE TRANSFORMATION GENERATION

After learning the embedding space of transformed feature sets, the embedding space should accurately describe and represent all transformed feature sets, including the unobserved best transformed feature set. Hence we utilize gradient ascent to identify the best embedding vector with the highest performance to decode.

Step 1: Gradient-ascent Optimization to Identify The Optimal Feature Set Embedding Another benefit of learning an embedding performance evaluator in the representation step is to make differentiable gradient optimization possible. Particularly, we extract the gradient from the evaluator toward the direction of maximizing feature set performance. Formally, the gradient-ascent is defined by : $\hat{E} = E + \eta \frac{\partial \omega_{\theta}}{\partial E}$, where ω_{θ} is the evaluator, \hat{E} is the optimal embedding, η is the step size. In the experiments, we select top-T transformed feature set embeddings in training data as the initialization seeds of gradient ascent. Therefore, we can identify T-improved embedding vectors as the optimal embedding set $\hat{\mathcal{E}} = \{\hat{E}^t\}_{t=1}^T$.

309 Step 2: Incorporating Flatness-aware into Gradient Ascent to Mitigate Shifts. In gradient op-310 timization, the flatness of the loss landscape has been proven to exhibit a close connection with 311 distribution shift resistance both theoretically and empirically. Under an open world, when there 312 exists a distribution shift between training and testing data, the optimal transformed feature set em-313 bedding on test samples does not coincide with the optimal transformed feature set embedding found 314 on training samples. Flatness ensures that the loss does not increase significantly in a neighborhood 315 around the found minimum. Therefore, flatness leads to distribution shift resistance because the loss on test examples does not increase significantly. Inspired by (Izmailov et al., 2018; Garipov 316 et al., 2018), we leverage the advantage of loss flatness and develop a flatness-aware gradient ascent 317 approach. 318

Specifically, unlike classic gradient ascent, we propose to enforce a searching process to oscillate around the optimal embedding to collect more suboptimal embedding vectors. This suboptimal embedding set can pinpoint a flat region where the real optimal point for the test set is located. We then aggregate all the suboptimal points as the final averaged embedding that represents the center of the flat region in the loss landscape. However, averaging all the suboptimal points of each gradient ascent iteration will introduce huge computational costs and include the non-optimal points, which 324 are far away from the optimal neighborhood. We therefore develop a cyclic scheme, in which a 325 cycle includes multiple gradient ascent iterations, and we only average the suboptimal points at the 326 end of a cycle. To avoid missing the most optimal point when approaching the maximum, we utilize 327 linearly decreasing learning rates within each cycle. In particular, Algorithm 3 shows we linearly 328 decrease the learning rate $\eta(i)$ from η_1 to η_2 over iterations during a cycle. In each iteration of a cycle, we initialize the learning rate and update the embedding by gradient ascent. At the end of 329 each cycle, we average the embedding (Line 7). 330

- 331
- Step 3: Decoding Embeddings to Reconstruct Optimal Feature Reconstruction Operations 332

After obtaining the candidate embedding set $\hat{\mathcal{E}}$, we use the well-trained decoder ψ to generate the 333 feature transformation operation sequences, i.e., $\hat{\mathcal{E}} \stackrel{\psi}{\to} \{\hat{\Upsilon}_i\}_{i=1}^T$. Specifically, the decoder iteratively 334 generates the next token of a feature transformation operation sequence, such as a feature token, an 335 operator token, or a segmentation token 'SEP' in an autoregressive manner until it produces an end 336 of sequence token 'EOS'. Finally, we divide the generated sequence into multiple segments using 337 the 'SEP' token and transform each segment into a feature, resulting in the optimal transformed 338 feature set with the best estimated performance. 339

INTEGRATING NORMALIZATION-DENORMALIZATION INTO PRE- AND 4.4 340 POST-PROCESSING 341

In addition to embedding dimension orthogonality and flatness-aware optimization, we propose a 342 data-centric perspective to alleviate distribution shift by first aligning and then recovering the distri-343 bution gaps between training and test data in pre and post-processing. We propose to incorporate a 344 normalization and denormalization mechanism in pre and post-processing. Specifically, we will first 345 normalize (e.g., z-score) original training and testing data before feature transformation, to obtain 346 normalized training and testing data with zero mean and one standard deviation. We then train our 347 method with the normalized training data to generate the optimal feature transformation operation 348 sequence, under a normalized distribution. Later, we apply the optimal feature transformation op-349 eration sequence learned from the normalized distribution to transform the normalized testing data. 350 So we obtain the transformed test data under the normalized distribution. Finally, we denormalize 351 (e.g., reverse z-score normalization) the transformed testing data to obtain the transformed testing 352 data in the original distribution. The underlying idea is to leverage normalization-denormalization 353 to conduct feature transformation learning in a normalized distribution and avoid shift-caused bias.

354 355

5 **EXPERIMENTAL RESULTS**

356 5.1 EXPERIMENTAL SETUP 357

Datasets. We perform experiments on 16 publicly available datasets from UCI (Dua & Graff, 2017) 358 and OpenML (Vanschoren et al., 2013). 9 datasets are intended for regression, and the remain-359 der parts are for classification. To evaluate the robustness of SAFT against distribution shifts, we 360 iteratively generate training and testing sets, employing the Kolmogorov-Smirnov test to identify 361 distribution shifts. Once a shift is detected, we finalize the training and testing sets with allocations 362 of 80% and 20%, respectively. More details about the data selection are included in Appendix C.1. 363

- Baselines. We compare our framework with 8 wildly-used feature transformation methods. We 364 described the details about all baselines in Appendix C.2.
- 366 Evaluation Metrics And Hyperparameters Setting. To control the variance of the downstream 367 model impact on evaluation, we apply Random Forests (RF) for both classification and regression 368 tasks. We described more details in Appendix C.3 and C.4.
- 369 Environmental Settings. All experiments are conducted on the Ubuntu 22.04.3 LTS operating 370 system, Intel(R) Core(TM) i9-13900KF CPU@ 3GHz, and 1 way RTX 4090 and 32GB of RAM, 371 with the framework of Python 3.11.4 and PyTorch 2.0.1. 372
- 5.2 PERFORMANCE RESULTS 373

374 In this experiment, we compare the feature transformation performance of SAFT with other baseline 375 models. Table 1 shows the overall comparison results in terms of F1 score and 1-RAE. We observe that SAFT consistently outperforms other baseline models across all datasets. The underlying driver 376 for this observation is that SAFT can compress the feature learning knowledge into a robust embed-377 ding space through shift-resistant embedding learning, and smoothly search for the optimal feature

Table 1: The overall performance results across various real-world datasets. The best and secondbest outcomes are indicated by bold and underlined fonts, respectively. We measure the performance
on classification (C) and regression (R) tasks using F1-score and (1-RAE) metrics, respectively. A
higher value indicates a superior quality of the feature transformation space.

-		_							-			
Dataset	C/R	Samples	Features	RDG	ERG	LDA	AFAT	NFS	TTG	GRFG	MOAT	SAFT
Housing Boston	R	506	13	0.375	0.366	0.146	0.387	0.395	0.383	0.361	0.395	0.405
Airfoil	R	1503	5	0.733	0.695	0.522	0.742	0.742	0.738	0.614	0.724	0.743
openml_586	R	1000	25	0.542	0.536	0.104	0.540	0.543	0.543	0.334	0.616	0.649
openml_589	R	1000	50	0.509	0.472	0.099	0.467	0.470	0.469	0.436	0.496	0.582
openml_607	R	1000	50	0.306	0.310	0.054	0.344	0.358	0.354	0.371	0.424	0.516
openml_616	R	500	50	0.197	0.311	0.014	0.342	0.344	0.343	0.459	0.369	0.534
openml_618	R	1000	50	0.421	0.391	0.058	0.436	0.430	0.431	0.257	0.427	0.468
openml_620	R	1000	25	0.510	0.464	0.025	0.475	0.464	0.462	0.495	0.566	0.545
openml_637	R	500	50	0.265	0.340	0.015	0.365	0.344	0.339	0.380	0.381	0.424
Higgs Boston	С	50000	28	0.693	0.694	0.508	0.692	0.691	0.696	0.698	0.702	0.704
SpectF	С	267	44	0.674	0.792	0.651	0.643	0.716	0.672	0.728	0.766	0.799
UCI Credit	С	30000	25	0.809	0.808	0.743	0.805	0.805	0.803	0.797	0.808	0.816
Wine Quality Red	С	999	12	0.658	0.491	0.588	0.662	0.650	0.675	0.668	0.681	0.700
Wine Quality White	С	4900	12	0.730	0.726	0.602	0.715	0.724	0.718	0.680	0.730	0.734
PimaIndian	С	768	8	0.679	0.761	0.685	0.636	0.691	0.734	0.689	0.763	0.780
Geman Credit	С	1000	24	0.697	0.662	0.693	0.656	0.698	0.683	0.647	0.707	0.743
BO	0.80 0.75	Recior Real F1-Score	Tr.m 0.75	icision Recall	SAFT SAFTw SAFTn SAFTn SAFT1	0.70 0.65 0.60 0.55	SAFT SAFT SAFT SAFT SAFT SAFT	0.8	1-MAE 1-MSE	SAFT 0. SAFT-w 0. SAFT-1 0. SAFT-1 0. 0. 1-RAE 0.	65 50 50 45 40 1-MAE 1-M	SAFT SAFT-W SAFT-W SAFT-M SAFT-I SAFT-I SAFT-I SAFT-I SAFT-I SAFT-I SAFT-I SAFT-W SAFT
(a) SpectF	(b)	Wine Whi	te (c)	Wine	Red	(d) Op	enml_61	5 (e)	Openm	1_618	(f) Open	ml_637

Figure 4: Results on the impact of normalization, Flatness Aware Gradient Ascent, and reweighting.

space via flatness-aware weight averaging. Moreover, we notice that the second-best model is variant among different cases. A potential reason for this observation is that baseline approaches overlook the impact of distribution shifts, resulting in an unstable identification of the optimal feature space. This experiment underscores that SAFT effectively captures invariant knowledge in feature learning against distribution shifts and produces a robust transformed feature space.

408 5.3 ABLATION STUDY

To evaluate the necessity of various components in SAFT, we develop three model variants: 1) SAFT-f, which excludes the Flatness Aware Gradient Ascent without performing weight averaging; 2) SAFT-n, which omits the normalization process; 3) SAFT-w, which removes the optimization of graph weights. We select three classification and three regression datasets for conducting ex-periments. Figures 4 show the comparison results. First, we observe that across all situations, the performance of the three model variants shows a decline when compared to SAFT. A potential rea-son is that disregarding strategies aware of distribution shifts fails to consider the variation between training and testing sets, resulting in suboptimal feature transformation performance. Moreover, SAFT can perform better compared to SAFT-w. This observation indicates that the learned graph weights can eliminate spurious correlations among features, thereby enhancing the transformed fea-ture space. Additionally, we find that SAFT outperforms SAFT-f. The underlying driver is that Flatness Aware Gradient Ascent prevents the gradient search process from converging to local opti-mal points, thereby smoothing and making the search process more robust. Furthermore, SAFT-n is inferior to SAFT. This observation reflects that aligning statistical properties via normalization is an effective strategy for tackling the OOD issue in feature space. Thus, these experiments demonstrate that each technical component of SAFT is indispensable for mitigating the impact of distribution shift on feature transformation.

5.4 ROBUSTNESS ANALYSIS

Dataset Split Robustness. To check the robustness of SAFT for different data splits, we use five
 distinct splits on the Wine Quality White dataset for evaluation. Each split defines a unique 20%
 proportion of the dataset as the testing set, moving sequentially through the dataset from beginning
 to end, with the remainder serving as the training set.

Figure 5 presents the comparison results in terms of F1-score, Precision, and Recall.



439 Figure 7: Comparison of the changes in joint and Figure 8: Comparison of the changes in joint and 440 in the Housing Boston dataset for both training in the Airfoil dataset for both training and test 442 443 After normalization.

We notice that despite potential distribution 445 shifts introduced by different dataset split ap-446 proaches, SAFT still demonstrates robust per-447 formance in addressing these shifts. The main 448 reason for this observation is that SAFT can 449 preserve the invariant and robust feature learn-450 ing knowledge within the embedding space for 451 robust and smooth search. Such a learning 452 mechanism can alleviate the impact of distri-453 bution shifts in different data splits. Hence, 454 this experiment shows that SAFT exhibits ro-455

456 Downstream Performance Robustness. To 457 check the robustness of SAFT for distinct 458 downstream ML models, we substitute the 459 downstream model with Random Forest (RF), 460 Support Vector Machine (SVM), K-Nearest 461 Neighborhood (KNN), Ridge, LASSO, and Decision Tree (DT). Table 2 shows the compari-462 son results on the Wine Quality Red dataset in 463 terms of the F1-score. We find that SAFT out-464 performs baselines in most cases. A potential 465 reason for this observation is that SAFT cap-466 tures invariant and generalizable feature learn-467 ing knowledge, leading to its superior general-468 ization across various downstream ML models. 469 Thus, this experiment demonstrates the robust-470 ness of SAFT across various ML models.

471

441

444

QUALITATIVE ANALYSIS 472 5.5

473 Shift Elimination. We perform a case study 474 to demonstrate the impact of normalization in



marginal distributions of age and medv (target) marginal distributions of delta and SSPL (target) and test sets. Left: Before normalization. Right: sets. Left: Before normalization. Right: After normalization.

Table 2: Robustness check on ML models for the Wine Quality White dataset.

	RF	SVM	KNN	DT	LASSO	Ridge
RDG	0.658	0.632	0.629	0.630	0.634	0.606
ERG	0.491	0.354	0.358	0.413	0.470	0.522
LDA	0.588	0.439	0.512	0.579	0.390	0.410
AFAT	0.662	0.582	0.443	0.663	0.596	0.587
NFS	0.650	0.634	0.542	0.639	0.634	0.601
TTG	0.675	0.644	0.525	0.602	0.630	0.620
GRFG	0.668	0.313	0.502	0.616	0.551	0.540
SAFT	0.700	0.655	0.556	0.677	0.659	0.658

bust performance in addressing distribution shift, regardless of the dataset split methods.



Figure 5: Robustness of the dataset split methods.



Figure 6: Comparison of the most important features in the original feature space of training and test set, and the SAFT generated feature space.

475 alleviating the distribution shift between the training and test sets. Figures 7 and 8 indicate the 476 changes in the distribution of the Housing Boston and Airfoil datasets, respectively. It is evident that, 477 after normalization, the joint and marginal distributions of age and medv (target) in both the training and test sets exhibit significant alignment. While the joint distribution of delta and SSPL (target) 478 remains relatively stable, the marginal distributions of the training and test sets show alignment. This 479 case study shows that normalization is effective in mitigating the distribution differences between 480 the training set and test set. 481

482 Feature Importance. We choose the top 10 most essential features from the original training set, 483 the original test set, and the transformed feature space using SAFT for the openml_616 dataset to facilitate comparison. The importance of features is measured by the mutual information (MI) be-484 tween the features and the target. Figure 6 displays the results. The labels accompanying each pie 485 indicate the respective feature names, with larger areas signifying greater importance. We observe

that the critical features in the training set differ from those in the test set, attributable to the distribution shift. Consequently, optimal feature combinations in the training set may not directly apply to
the test set, as their significance may vary. Furthermore, we note that in the new feature space generated by SAFT only 7 features are employed (with 2 new features generated by SAFT), but these
features improve the performance of ML model by 29.24%. This phenomenon illustrates that SAFT
can identify and generate effective features that are applicable to the test set. SAFT can capture the
true representation of the feature set and generate a robust embedding space.

To make the experiment more convincing, we also analyzed the time complexity and space complexity in Appendix D. Meanwhile, we visualize the embedding space to verify the effectiveness of the intermediate encoding in Appenddix E.

496 497

6 RELATED WORKS

498 Feature Transformation aims to construct an enhanced feature space by transforming original fea-499 tures. Previous research can be categorized into three categories: 1) Expansion-reduction based 500 approaches (Kanter & Veeramachaneni, 2015; Horn et al., 2020; Khurana et al., 2016), where the 501 original feature space is expanded through explicit or greedy mathematical transformations and sub-502 sequently the feature space is narrowed down by selecting valuable features. 2) Evolution-evaluation 503 approaches (Wang et al., 2022; Khurana et al., 2018; Tran et al., 2016), where evolutionary algo-504 rithms or Reinforcement learning models are employed to optimize the process of iteratively creating 505 effective features and retaining significant ones. 3) Auto ML-based approaches (Chen et al., 2019; Zhu et al., 2022; Wang et al., 2023), where the most appropriate model architecture is automatically 506 identified to formulate AFT as an Auto ML task. However, these methods face two challenges: 1) 507 high-order feature transformation is difficult to produce; 2) transformation performance is unsta-508 ble. To address these obstacles, we formulate AFT as a continuous optimization task (Wang et al., 509 2023). In particular, we employ an RL-based data collector to prepare the training data and utilize 510 an encoder-decoder evaluator-based architecture to build the continuous space. Then, we search for 511 the optimum solution within the continuous space and reconstruct the transformed feature space. 512

Distribution Shift Problems, arising from inconsistent distributions between the training set and 513 test set, is prevalent in the real-world domains. Existing works to solve this problem mainly focus on 514 time series analysis (Fan et al., 2023), computer vision (Yu et al., 2023), and natural language pro-515 cessing (Dou et al., 2022), etc. But distribution shift problem for feature transformation is relatively 516 underexplored and can greatly degrade the effectiveness of the reconstructed feature. To tackle this 517 challenge, specifically, we integrate three techniques designed to function at different levels within 518 our framework: 1) shift-resistant representation: where the invariant true representation is learned 519 under DSFT; 2) flatness-aware generation: where a more robust embedding is pinpointed within the 520 searching process; and 3) shift-aligned pre and post-processing: where distributions are aligned for 521 the data processing.

7 CONCLUSION

524 In this work, we combine representation, generation, and anti-shift learning to design a robust feature transformation framework for DSFT, namely Shift Awareness Feature Transformation (SAFT) 526 framework. Our contributions can be summarized as follows: First, to avoid searching in a large 527 discrete, we formulate DSFT as a gradient optimization problem and develop a representation-528 generation framework to identify optimal transformed features. Second, to eliminate the effect of 529 distribution shift, we integrate three mechanisms designed to operate at varying levels within our 530 framework: 1) shift-resistant representation: to learn the invariant true representation under DSFT; 2) flatness-aware generation: to identify a more robust embedding in the searching process; and 3) 531 shift-aligned pre and post processing: to align distributions for the data processing. Ultimately, the 532 experimental results empirically demonstrate that even without the prior knowledge of test distribu-533 tion, SAFT can generate a robust feature space and identify optimal transformed feature sequences 534 based on the information learned from the training set. This emphasizes the superior generalization 535 ability of SAFT to enhance AI capabilities in diverse open environments across domains, such as 536 economics and health care. 537

538

522 523

520

540 REFERENCES

547

578

579

580

- Susan Athey, Guido W Imbens, and Stefan Wager. Approximate residual balancing: debiased inference of average treatment effects in high dimensions. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 80(4):597–623, 2018.
- Vance W Berger and YanYan Zhou. Kolmogorov–smirnov test: Overview. Wiley statsref: Statistics
 reference online, 2014.
- David M Blei, Andrew Y Ng, and Michael I Jordan. Latent dirichlet allocation. *Journal of machine Learning research*, 3(Jan):993–1022, 2003.
- Xiangning Chen, Qingwei Lin, Chuan Luo, Xudong Li, Hongyu Zhang, Yong Xu, Yingnong Dang,
 Kaixin Sui, Xu Zhang, Bo Qiao, et al. Neural feature search: A neural architecture for automated
 feature engineering. In 2019 IEEE International Conference on Data Mining (ICDM), pp. 71–80.
 IEEE, 2019.
- Shihan Dou, Rui Zheng, Ting Wu, Songyang Gao, Junjie Shan, Qi Zhang, Yueming Wu, and Xuan jing Huang. Decorrelate irrelevant, purify relevant: Overcome textual spurious correlations from a feature perspective. *arXiv preprint arXiv:2202.08048*, 2022.
- 558 Dheeru Dua and Casey Graff. Uci machine learning repository, 2017. URL http://archive. ics.uci.edu/ml.
- Wei Fan, Pengyang Wang, Dongkun Wang, Dongjie Wang, Yuanchun Zhou, and Yanjie Fu. Dish-ts:
 a general paradigm for alleviating distribution shift in time series forecasting. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pp. 7522–7529, 2023.
- Wei Fan, Shun Zheng, Pengyang Wang, Rui Xie, Jiang Bian, and Yanjie Fu. Addressing distribution shift in time series forecasting with instance normalization flows. *arXiv preprint arXiv:2401.16777*, 2024.
- Timur Garipov, Pavel Izmailov, Dmitrii Podoprikhin, Dmitry P Vetrov, and Andrew G Wilson. Loss
 surfaces, mode connectivity, and fast ensembling of dnns. *Advances in neural information processing systems*, 31, 2018.
- 570 Cyril Goutte and Eric Gaussier. A probabilistic interpretation of precision, recall and f-score, with implication for evaluation. In *European conference on information retrieval*, pp. 345–359.
 572 Springer, 2005.
- Will Hamilton, Zhitao Ying, and Jure Leskovec. Inductive representation learning on large graphs.
 Advances in neural information processing systems, 30, 2017.
- Arthur V Hill. The encyclopedia of operations management: a field manual and glossary of operations management terms and concepts. Ft Press, 2012.
 - Timothy O Hodson. Root-mean-square error (rmse) or mean absolute error (mae): When to use them or not. *Geoscientific Model Development*, 15(14):5481–5487, 2022.
- Franziska Horn, Robert Pack, and Michael Rieger. The autofeat python library for automated feature engineering and selection. In *Machine Learning and Knowledge Discovery in Databases: International Workshops of ECML PKDD 2019, Würzburg, Germany, September 16–20, 2019, Proceedings, Part I*, pp. 111–120. Springer, 2020.
- Xuanming Hu, Wei Fan, Kun Yi, Pengfei Wang, Yuanbo Xu, Yanjie Fu, and Pengyang Wang. Boosting urban prediction via addressing spatial-temporal distribution shift. In 2023 IEEE International Conference on Data Mining (ICDM), pp. 160–169. IEEE Computer Society, 2023.
- Pavel Izmailov, Dmitrii Podoprikhin, Timur Garipov, Dmitry Vetrov, and Andrew Gordon Wilson. Averaging weights leads to wider optima and better generalization. *arXiv preprint arXiv:1803.05407*, 2018.
- James Max Kanter and Kalyan Veeramachaneni. Deep feature synthesis: Towards automating data science endeavors. In 2015 IEEE international conference on data science and advanced analytics (DSAA), pp. 1–10. IEEE, 2015.

594 Udayan Khurana, Deepak Turaga, Horst Samulowitz, and Srinivasan Parthasrathy. Cognito: Auto-595 mated feature engineering for supervised learning. In 2016 IEEE 16th International Conference 596 on Data Mining Workshops (ICDMW), pp. 1304–1307. IEEE, 2016. 597 Udayan Khurana, Horst Samulowitz, and Deepak Turaga. Feature engineering for predictive mod-598 eling using reinforcement learning. In Proceedings of the AAAI Conference on Artificial Intelligence, volume 32, 2018. 600 601 Taesung Kim, Jinhee Kim, Yunwon Tae, Cheonbok Park, Jang-Ho Choi, and Jaegul Choo. Reversible instance normalization for accurate time-series forecasting against distribution shift. In 602 International Conference on Learning Representations, 2021. 603 604 Andrew Kusiak. Feature transformation methods in data mining. IEEE Transactions on Electronics 605 packaging manufacturing, 24(3):214–221, 2001. 606 Haoyang Li, Xin Wang, Ziwei Zhang, and Wenwu Zhu. Ood-gnn: Out-of-distribution generalized 607 graph neural network. IEEE Transactions on Knowledge and Data Engineering, 2022. 608 609 Xingyu Li, Bo Tang, and Haifeng Li. Adaer: An adaptive experience replay approach for continual 610 lifelong learning. Neurocomputing, 572:127204, 2024. 611 Andrzej Maćkiewicz and Waldemar Ratajczak. Principal components analysis (pca). Computers & 612 Geosciences, 19(3):303-342, 1993. 613 614 Zheyan Shen, Peng Cui, Tong Zhang, and Kun Kunag. Stable learning via sample reweighting. In 615 Proceedings of the AAAI Conference on Artificial Intelligence, volume 34, pp. 5692–5699, 2020. 616 Binh Tran, Bing Xue, and Mengjie Zhang. Genetic programming for feature construction and se-617 lection in classification on high-dimensional data. *Memetic Computing*, 8:3–15, 2016. 618 619 Joaquin Vanschoren, Jan N. van Rijn, Bernd Bischl, and Luis Torgo. Openml: networked science in machine learning. SIGKDD Explorations, 15(2):49-60, 2013. doi: 10.1145/2641190.2641198. 620 URL http://doi.acm.org/10.1145/2641190.264119. 621 622 Dongjie Wang, Yanjie Fu, Kunpeng Liu, Xiaolin Li, and Yan Solihin. Group-wise reinforcement fea-623 ture generation for optimal and explainable representation space reconstruction. In Proceedings of 624 the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, pp. 1826–1834, 625 2022. 626 Dongjie Wang, Meng Xiao, Min Wu, Pengfei Wang, Yuanchun Zhou, and Yanjie Fu. 627 Reinforcement-enhanced autoregressive feature transformation: Gradient-steered search in con-628 tinuous space for postfix expressions. arXiv preprint arXiv:2309.13618, 2023. 629 Dongjie Wang, Meng Xiao, Min Wu, Yuanchun Zhou, Yanjie Fu, et al. Reinforcement-enhanced 630 autoregressive feature transformation: Gradient-steered search in continuous space for postfix 631 expressions. Advances in Neural Information Processing Systems, 36, 2024. 632 633 Wangyang Ying, Dongjie Wang, Kunpeng Liu, Leilei Sun, and Yanjie Fu. Self-optimizing feature 634 generation via categorical hashing representation and hierarchical reinforcement crossing. arXiv 635 preprint arXiv:2309.04612, 2023. 636 Runpeng Yu, Songhua Liu, Xingyi Yang, and Xinchao Wang. Distribution shift inversion for out-637 of-distribution prediction. In Proceedings of the IEEE/CVF Conference on Computer Vision and 638 Pattern Recognition, pp. 3592-3602, 2023. 639 640 Daochen Zha, Zaid Pervaiz Bhat, Kwei-Herng Lai, Fan Yang, Zhimeng Jiang, Shaochen Zhong, and 641 Xia Hu. Data-centric artificial intelligence: A survey. arXiv preprint arXiv:2303.10158, 2023. 642 Peiyan Zhang, Yuchen Yan, Chaozhuo Li, Senzhang Wang, Xing Xie, Guojie Song, and 643 Sunghun Kim. Continual learning on dynamic graphs via parameter isolation. arXiv preprint 644 arXiv:2305.13825, 2023. 645 Xingxuan Zhang, Peng Cui, Renzhe Xu, Linjun Zhou, Yue He, and Zheyan Shen. Deep stable 646 learning for out-of-distribution generalization. In Proceedings of the IEEE/CVF Conference on 647 Computer Vision and Pattern Recognition, pp. 5372-5382, 2021.

648	Guanghui Zhu Zhuaar Yu Chunfang Yuan and Yihua Huang Difar: differentiable automated
649	feature engineering. In International Conference on Automated Machine Learning, pp. 17–1.
650	PMLR, 2022.
651	
652	
653	
654	
655	
656	
657	
658	
659	
660	
661	
662	
663	
664	
665	
666	
667	
668	
669	
670	
671	
672	
673	
674	
675	
676	
677	
678	
679	
680	
681	
682	
683	
684	
685	
686	
687	
688	
689	
690	
691	
692	
693	
694	
695	
696	
697	
698	
699	
700	
701	

A MORE DETAILS ABOUT DATA COLLECTION BY RL AGENTS

In this section, we provide more details on the RL agents used for data collection.

705 Leveraging reinforcement learning to explore high-quality and diverse training data. The train-706 ing data collector includes: 1) Multiple Agents: We design three agents to perform feature crossing: 707 a head feature agent, an operation agent, and a tail feature agent. 2) Actions: In each reinforcement 708 iteration, the three agents collaborate to select a head feature, an operator (e.g., +,-,*,/), and a tail 709 feature to generate a new feature (e.g. f_1+f_2). The newly generated feature is later added to the 710 feature set for the next feature generation. To balance diversity and quality, we employ the ϵ -greedy 711 DQN algorithm. For each step, each agent has a ϵ probability of selecting a feature/operation based 712 on policy and has a $(1-\epsilon)$ probability of selecting a feature/operation equally. The ϵ will increase 713 over time so that agents can first explore more and then exploit more. 3) Environment: The en-714 vironment is the feature space, representing an updated feature set. When three feature agents 715 generate and add a new feature to the previous feature set, the state of the feature space (i.e., the 716 environment) changes. The state represents the statistical characteristics (such as aggregated average mean, and variance) of the feature subspace. 4) Reward Function: We formulate the reward as 717 the improvement of the performance (e.g., accuracy) of the explored feature set on a downstream 718 task (e.g., regression) on the current iteration, compared with that of the explored feature set on 719 the previous iteration. In this way, we incentivize agents to explore high-quality feature sets. 5) 720 Training and Optimization: Our reinforcement data collector includes many feature crossing steps. 721 Each step consists of two stages - control and training. In the control stage, each feature agent takes 722 actions to change the size and contents of a new feature set and generate a reward. This reward is 723 assigned to each participating agent. In the training stage, the agents train their policies via experi-724 ence replay independently by minimizing the mean squared error (MSE) of the Bellman Equation. 725 6) Feature Set-Performance Annotation Pairs Preparation: We test each RL-explored feature subset 726 with a downstream ML task (e.g., random forest classification) to collect feature set performance 727 annotation. We select the top 5,000 feature set-performance annotation pairs based on predictive accuracy as training data. 728

729 Training data triples. After obtaining the transformed feature set-performance pairs through RL, 730 we convert this data into triples to build our training dataset. 1) We convert the transformed feature 731 set into a feature-feature graph as input to the encoder. The purpose is to capture the structural 732 relationships between features, enabling the encoder to better represent complex feature interactions. 733 2) The transformed feature set is treated as a feature transformation operation sequence. The purpose of obtaining this sequence is to allow the decoder to recognize and process it. The decoder then 734 iteratively predicts one token at a time until it reaches an end-of-sequence (EOS) token. 3) The 735 evaluator assesses the performance associated with the sequence, which guides the gradient-based 736 search to identify the optimal embedding. 737

The determination of the length of the feature transformation operation sequence. We introduce a hyperparameter to control the number of transformed features. When this number exceeds the predefined limit, feature selection is applied to retain only the most important features. At each step, RL generates a new crossed feature, which is added to the original feature set as a candidate for the next iteration. Thus, the length of the sequence is influenced by the order in which the transformed features are added to the set. For example, two transformed feature sets have the same number of features but different lengths of sequence (e.g., $SOSf_1SEPf_2 + f_1EOS$ vs $SOSf_1SEP\sqrt{f_2}EOS$).

745 746

747

748 749

750 751

752 753

755

704

B MORE DETAILS ABOUT SHIFT AWARENESS FEATURE TRANSFORMATION FRAMEWORK (SAFT)

B.1 ALGORITHM TABLE FOR THE OVERALL FRAMEWORK

This part provides an algorithm table for the overall framework (see Algorithm 1).

754 B.2 Algorithm Table for the Shift-resistant Bilevel Training

In this part, we provide an algorithm table for the shift-resistant bilevel training (see Algorithm 2).



⁰ Algo	prithm 2 Shift-resistant Bilevel Training
$1 \frac{v}{Inv}$	ut: EPOCH_1 and EPOCH_2
2 Out	put: Learned Embedding Space
³ 1: 1	for out_loop_epoch $\leftarrow 1$ to EPOCH_1 do
4 2:	Forward propagate
5 3:	for inner_loop_epoch $\leftarrow 1$ to EPOCH_2 do
6 4:	Optimize the sample weighting by Equation (1)
7 5:	end for
B 6:	Back propagate with the weighted joint loss $\mathcal{L} = \mathcal{L}_{est} + \gamma \mathcal{L}_{rec}$ to joint train
7: 0	end for
0	
	withm 2 Elethors Awara Gradient Assent
$2 \frac{\text{Alge}}{T}$	
Inp	ut: Initialized embedding E, LR bounds η_1, η_2 , cycle length c, number of iterations n
	put: Optimal embedding E
1:	$E \leftarrow E$
2: 1	for $i = 1, 2,, n$ do
3:	$\eta \leftarrow \eta(i) \ \{ \text{ Calculate LR for the iteration, } \eta_1 < \eta(i) < \eta_2 \} $
4:	$E \leftarrow E + \eta \frac{\partial \omega_{\theta}}{\partial E} \{ \text{Stochastic gradient update} \}$
5:	if mod $(i, c) = 0$ then
6:	$n_{\text{models}} \leftarrow \frac{i}{c} \{ \text{Number of models} \}$
7:	$\hat{E} \leftarrow \frac{\hat{E} \cdot n_{\text{models}} + E}{4} \{ \text{Update average} \}$
	$n_{\rm models}+1$ (r
×٠	end If

we use five distinct splits on the Wine Quality White dataset to evaluate the proposed framework to respond to varying degrees of distributional shifts.

C.2 MORE DETAILS FOR THE BASELINES

We compare our framework with 8 wildly-used feature transformation methods: 1) Random Gener-841 ation (RDG) randomly produces feature-operation-feature transformations to generate a new feature 842 space; 2) Essential Random Generation (ERG) initially expands the feature space by applying op-843 erations to each feature, and then chooses the essential features as the new feature space; 3) Latent 844 Dirichlet Allocation (LDA) (Blei et al., 2003) refines the feature space to get the factorized hidden 845 state via matrix factorization; 4) AutoFeat Automated Transformation (AFAT) (Horn et al., 2020) 846 is an enhanced version of ERG, which repeatedly generates new features and employs multi-step 847 feature selection to identify informative ones; 5) Neural Feature Search (NFS) (Chen et al., 2019) 848 generates the transformation sequence for each feature and the entire process is optimized by RL; 849 6) Traversal Transformation Graph (TTG) (Khurana et al., 2018) formulates the transformation 850 process as a graph and subsequently employs an RL-based search method to find the optimal feature set; 7) Group-wise Reinforcement Feature Generation (GRFG) (Wang et al., 2022) employs 851 three collaborative reinforced agents to perform feature generation for feature space refinement. 8) 852 reinforceMent-enhanced autOregressive feAture Transformation (MOTA) (Wang et al., 2024) 853 formulates the discrete feature transformation problem as a continuous optimization task, while still 854 assuming the I.I.D. condition. 855

856 857

858

833 834 835

836

837 838 839

840

C.3 EVALUATION METRICS

To control the variance of the downstream model impact on evaluation, we apply Random Forests
(RF) for both classification and regression tasks. For classification, we use the F1-score, Precision, and Recall as the evaluation metrics (Goutte & Gaussier, 2005). In regression, we assess
performance using the following metrics (Hill, 2012; Hodson, 2022) 1 - Relative Absolute Error (1-RAE), 1 - Mean Absolute Error (1-MAE), and 1 - Mean Squared Error (1-MSE). For all these metrics, a higher value indicates a more effective feature transformation space.

Table 5: Time cost comparisons with baseline	Table 3:	le 3: Time cos	t comparisons	with	baselines
--	----------	----------------	---------------	------	-----------

Dataset	RL Data Collection (min)	RDG(s)	LDA(s)	ERG(s)	AFAT(s)	NFS(s)	TTG(s)	GRFG(s)	MOTA(min)	SAFT(min)
SpectF	33.7	3.0	0.3	9.5	3.8	8.3	9.1	282.1	156.0	142.1
Wine Quality White	114.7	30.8	0.5	29.4	4.2	24.7	26.4	373.2	195.8	39.9
Wine Quality Red	36.5	7.4	0.6	13.9	3.6	9.2	10.3	182.2	107.5	43.1
openml_616	156.6	27.5	0.2	55.9	2.6	18.1	24.6	936.5	430.1	160.7
openml_618	239.6	67.7	0.2	82.4	3.0	30.9	38.7	821.6	373.8	95.9
openml_637	171.3	28.1	0.2	55.8	2.8	19.3	25.1	731.7	300.9	147.9

Table 4: Space complexity comparisons with baselines.

Dataset	RL Data Collection (MB)	RDG(MB)	LDA(MB)	ERG(MB)	AFAT(MB)	NFS(MB)	TTG(MB)	GRFG(MB)	MOTA(MB)	SAFT(MB)
SpectF	0.19	0.13	0.07	0.14	0.08	0.14	0.13	1.04	0.14	0.16
Wine Quality White	0.21	0.15	0.08	0.16	0.09	0.15	0.14	18.70	0.13	0.48
Wine Quality Red	0.19	0.14	0.07	0.15	0.08	0.15	0.15	6.14	0.13	0.19
openml_616	0.20	0.14	0.07	0.16	0.08	0.16	0.14	1.94	0.14	0.17
openml_618	0.23	0.15	0.09	0.16	0.11	0.17	0.14	3.84	0.14	0.23
openml_637	0.23	0.16	0.10	0.17	0.11	0.17	0.15	1.94	0.14	0.17

C.4 HYPERPARAMETERS SETTING

1) <u>RL collector</u>: We use the reinforcement data collector to collect explored feature set-feature utility score pairs. The collector explores 512 episodes, and each episode includes 10 steps. 2) <u>Graph Construction</u>: The threshold for creating an edge between two features is the 95th percentile of all similarity values. 3) <u>Our framework</u>: we map the attribute of each node to a 64-dimensional embedding, and use a 2-layer GNN network and a 2-layer projection head to integrate such information. The decoder is a 1-layer LSTM network, which reconstructs a feature cross sequence. The evaluator is a 2-layer feed-forward network, in which the dimension of each layer is 200. During optimization, During optimization, we assign α (i.e., 10) as the weight for estimation loss and β (i.e., 0.1) as the weight for reconstruction loss, the batch size is 256, the epochs are 500, and the learning rate range is 0.001-0.0005.

D COMPARISON OF TIME AND SPACE COMPLEXITIES WITH BASELINE MODELS.

Time complexity. We analyzed the computational complexity of SAFT compared to baseline meth-ods across six datasets. Table 3 shows that while RDG, LDA, ERG, AFAT, NSF, and TTG require less time, their accuracy is lower due to their inability to handle the shift problem. In contrast, our method offers lower costs and better performance compared to GRFG and MOTA. Additionally: 1) Feature transformation is neither a critical step in data preparation/processing nor a major time factor for most tasks; 2) The time cost of our method is acceptable, especially when compared to manual feature engineering, which can take days or months. Our approach significantly saves time; 3) Although our method may take more time than some baselines, it achieves superior feature en-gineering performance; 4) With our encoder-optimization-generation (EOG) design, we can further reduce training time by pre-training a foundational model and then fine-tuning it to other domains.

We also report the computational overhead of RL data collection, as shown in Table 3. While SAFT
incurs additional time costs during data collection, this process is done asynchronously and offline,
with an acceptable time cost. The main reason is that the RL-based collector requires more time to
gather high-quality data, and the sequence formulas across the feature space increase the learning
time of the sequence model.

913 Space complexity. To analyze the space complexity of SAFT, we illustrate the storage size required
914 for different datasets. Table 4 shows that SAFT occupies minimal space and remains relatively
915 stable. This is primarily due to the EOG framework, which embeds knowledge from variable916 length discrete sequences into fixed-length embedding vectors. This embedding process keeps the
917 parameter size stable, preventing it from increasing with data growth. Therefore, the experiments demonstrate that SAFT exhibits good scalability across datasets of different sizes.



Figure 10: Visualization of the convergence embedding space. Each point represents a feature transformation operation sequence. The orange points represent the top 50 embedding points based on their downstream performance.

E VISUALIZATION FOR THE LEARNED EMBEDDING SPACE.

In this section, we visualized the embedding space to validate the effectiveness of the intermedi-ate encoding. First, we collected the latent embeddings generated from the transformation records. Then, we used t-SNE to project these embeddings into a two-dimensional space for visualization. We selected four datasets for visualization as shown in Figure 10, where each point represents a fea-ture transformation operation sequence. These points exhibit different distribution patterns, likely due to variations in the lengths of the corresponding transformation sequences, causing them to spread out in the embedding space. However, despite their differing positions, we observed that the points tend to cluster into several distinct groups, especially the points representing top-performing sequences (orange points), which tend to cluster closely together. This suggests that data points with strong downstream task performance are likely concentrated in certain regions, allowing gradient-based search algorithms to effectively locate optimal points within these areas when we use these points as initial search seeds. This case study highlights the role of continuous space in extract-ing feature knowledge through reconstruction and loss estimation, enabling the search for optimal embeddings to reconstruct the best feature transformation sequences.