LLM-GUIDED SELF-SUPERVISED TABULAR LEARN ING WITH TASK-SPECIFIC PRE-TEXT TASKS

Anonymous authors

Paper under double-blind review

Abstract

One of the most common approaches for self-supervised representation learning is defining pre-text tasks to learn data representations. Existing works determine pre-text tasks in a "task-agnostic" way, without considering the forthcoming downstream tasks. This offers an advantage of broad applicability across tasks, but can also lead to a mismatch between task objectives, potentially degrading performance on downstream tasks. In this paper, we introduce TST-LLM, a framework that effectively reduces this mismatch when the natural language-based description of the downstream task is given without any ground-truth labels. TST-LLM instructs the LLM to use the downstream task's description and meta-information of data to discover features relevant to the target task. These discovered features are then treated as ground-truth labels to define "target-specific" pre-text tasks. TST-LLM consistently outperforms contemporary baselines, such as STUNT and LFR, with win ratios of 95% and 81%, when applied to 22 benchmark tabular datasets, including binary and multi-class classification, and regression tasks.

1 INTRODUCTION

027 028

004

010 011

012

013

014

015

016

017

018

019

021

024 025 026

029 Obtaining unlabeled data for machine learning is typically more scalable and cheaper than gathering labeled data in real-world applications. Self-supervised representation learning, which was proposed to extract useful information from unlabeled data, enhances the performance of downstream tasks by 031 obtaining superior representations (Chen et al., 2020; Oord et al., 2018; Tschannen et al., 2019). A common approach for self-supervised representation learning is utilizing a pre-text task based on 033 the type or characteristics of the data, and then learning representations by optimizing the objective 034 of the pre-text task (instead of the downstream task) (Assran et al., 2022; Kim et al., 2018; Zhang et al., 2017). For example, in computer vision domain, such a pre-text task can be defined to estimate the degree of rotation applied to an image (Gidaris et al., 2018). Other works have defined 037 augmentations that do not deform the contents of images, such as horizontal flips or color jittering, 038 to learn representations that are invariant to these modifications (Grill et al., 2020; Han et al., 2020; Zbontar et al., 2021).

These pre-text task-based methods have also been extended to tabular data – specifically, efforts have been made to adapt pre-text tasks that were previously limited to images and text. Common tasks for tabular data include corrupting or masking data and then reconstructing the original sample from the representation (Yoon et al., 2020; Wu et al., 2024), or designing augmentations suited to tabular data to perform contrastive learning tasks (Bahri et al., 2022; Somepalli et al., 2022). The inductive bias provided by the pre-text tasks plays a role in preemptively removing spurious correlations or noisy information within tabular data.

Despite the success of using pre-text tasks in tabular learning, a fundamental limitation remains in the task mismatch between the pre-text tasks optimized by representation learning and the actual downstream tasks (Sui et al., 2024). This is because the pre-text tasks in representation learning are determined in a "task-agnostic" way that are oblivious to the forthcoming downstream tasks.
Task-agnostic definitions offer the advantage of broad applicability across tasks, but also potentially risks including noisy irrelevant information or eliminating critical information for the downstream task. For instance, using additive Gaussian noise or affine transformations as augmentations in tabular data with high variable correlation can create unrealistic samples and simultaneously lose

ost
 correlation information (Sui et al., 2024; Hajiramezanali et al., 2022). Such fallacies impair the performance of the downstream task.

Our study seeks to address the issue of the task-objective mismatch by defining pre-text tasks in 057 a task-specific rather than task-agnostic manner using the natural-language based description of the downstream task. The description includes the task objective (e.g., "Does this person earn more than 50,000 dollars per year?") and the answer candidates (e.g., "Yes" or "No"). Using this information, 060 we propose, TST-LLM (Task-specific Self-supervised Tabular learning with LLMs), which aims to 061 improve representation learning via LLM-discovered features without incorporating any ground-truth 062 labels or label statistics. By leveraging the prior knowledge of the LLM, we explore the relationship 063 between the task and data features from their natural language-based descriptions. This process aims 064 to determine which combinations or transformations of original data features can yield meaningful information for solving the task. Then, the new features created through the LLM's prior knowledge 065 are used to generate the ground-truth labels for the pre-text tasks that train the representation. For 066 example, in the task of predicting whether a person earns more than 50,000 dollars, newly discovered 067 features such as "age * working hours", based on prior knowledge, are likely to have a 068 higher correlation with the label. Learning with these features provides additional task-relevant 069 information, such as the importance of original features to the task and the correlation between them.

071 TST-LLM consists of two main stages. In the first stage, target task's textual description, meta-information of data (e.g., feature names and descriptions), and text-serialized unlabeled data 072 are used to construct prompts. Then, they are fed into an LLM to extract new task-relevant features. 073 This process is repeated, while previously extracted features are excluded at the next iteration to 074 ensure the diversity of feature synthesis. In the second stage, the discovered features are considered 075 as ground-truth labels to define pre-text tasks. We use supervised contrastive learning (Khosla et al., 076 2020) to perform multi-task learning for each label, learning useful representations. Additionally, we 077 introduce a process for selecting a diverse feature set from the discovered features that are distinctly 078 aligned with both the original data and each other for computational efficiency. 079

A key advantage of TST-LLM is its simplicity; it can be applied to any problem as long as a task description and feature descriptions are defined in natural language. We demonstrate that the features discovered by the LLM are meaningful and relevant to the actual target task. Our model consistently outperforms contemporary baselines, such as STUNT and LFR, with win ratios of 95% and 81%, when applied to 22 benchmark tabular datasets including binary and multi-class classification, and regression tasks.

085 086 087

880

2 RELATED WORK

Self-supervised representation learning for tabular data. New advancements in self-supervised representation learning enable the discovery of meaningful representations from unlabeled data 091 across wide modalities, from images (Caron et al., 2020; Wen et al., 2022; Wu et al., 2018) to 092 texts (Gao et al., 2021; Kenton & Toutanova, 2019; Radford et al., 2019), audio (Mittal et al., 2022; Owens & Efros, 2018), and most recently to tabular data (Balestriero et al., 2023; Gharibshah & 094 Zhu, 2022). One of the common ways of self-supervised learning is to define a pre-text task on 095 an unlabeled dataset to facilitate learning. According to the literature (Gharibshah & Zhu, 2022), 096 pre-text tasks for tabular data can be broadly classified into three types. The first category, invariance 097 learning, involves defining a positive view of a given sample and learning the representation invariance 098 between them. The positive view of the sample can be created using weak augmentations that do not distort the original content (Bahri et al., 2022; Somepalli et al., 2022) or by selecting samples with similar characteristics from the training data (Nam et al., 2023b). The second category, predictive 100 learning, includes methodologies that generate explicit labels from the dataset and train the model 101 to predict these labels. For example, masking or corrupting data and then using the original data for 102 reconstruction as a label (Wu et al., 2024; Yoon et al., 2020). Some studies also proposed pre-text tasks 103 on various publicly available benchmark datasets or synthetic datasets and then performing transfer 104 learning for downstream tasks (Hollmann et al., 2023; Wang & Sun, 2022). The last category includes 105 a hybrid approach combining invariance and predictive learning (Ucar et al., 2021; Zhu et al., 2023). 106

107 All of the above methods define pre-text tasks in a task-agnostic manner, which can lead to inconsistencies with the actual objectives of the downstream tasks that can ultimately hinder performance.



Problem formulation. An unlabeled tabular dataset with *d*-dimensional input features $\mathcal{D} = \{\mathbf{x}_i\}_{i=1}^N$ is given where $\mathbf{x}_i \in \mathbb{R}^d$. A downstream task description in natural language, E_{task} , and the names and short descriptions of each feature, $E_{\text{name}} = \{e_{\text{name}}^j\}_{j=1}^d$ and $E_{\text{desc}} = \{e_{\text{desc}}^j\}_{j=1}^d$, are also provided. The model aims to train an encoder f that extracts informative data representations to tackle the downstream task in an unsupervised setting, i.e., no ground-truth labels are provided. The downstream tasks can be binary or multi-class classification and regression.

Figure 1 illustrates the entire process of our model. TST-LLM utilizes downstream task description and meta information to define pre-text tasks that aligns with the downstream task objectives, and performs representation learning via these tasks. Initially, the model passes the task description E_{task} along with meta-information of data E_{name} and E_{desc} to the LLM to generate potentially relevant features through combinations or transformations of original features (Section 3.1). These generated features are set as target labels for the pre-text tasks, which are then used to train the encoder through multi-task contrastive learning (Section 3.2). Details of each stage are described below.

158 159

3.1 LLM-GUIDED FEATURE DISCOVERY WITH TASK DESCRIPTION

161 The model generates data-related features from task description and meta information utilizing the prior knowledge and reasoning abilities of an LLM. Feature discovery process involves running

You are a data engineer. Given the task description and the list of features and data examples, you are
making a new column for the data which is informative to solve the task.
Task: [Downstream Task Description]
Features: [Feature Descriptions]
Examples: [Serialized Examples]
Given a type of operations below, generate 5 new columns which are the most informative to solve the
task using operations. Refer to the examples when generating features. Only use features listed in the
reature description. Note that multiple operations can be nested to generate a new column.
The possible type of operations is as follows:
[Operation Descriptions]
You also have some new example features generated with these modules.
Example Features:
Index realure_indine realure_desc [Features From Previous Trial]
You must write new feature that is different from all above examples features with respect to both
names and descriptions.
Format of response for 5 new columns:
Thought 1: Any reasons why the following new feature would be helpful for the task
New feature 1: Type of operation New_column_name One line pseudo code for generating columns
Thought 5:
New feature 5:

Figure 2: Prompt for feature discovery. Text in blue corresponds to data description summary part, red text to operation instruction, teal text to diversity enforcement, and brown text to response instruction part.

192 193 194

195

196

197

201

204 205

207

208

209

210

211 212

213

214

215

189

190

191

multiple LLM inferences. The designed prompt consists of four main components: data description summary, operation instruction, diversity enforcement, and response instruction (see Figure 2 and Appendix A.2 for the example prompt).

Data description summary. This component provides a basic data description for feature discovery 198 (see blue part in Figure 2). It includes the downstream task's description E_{task} (e.g., "Does this person 199 earn more than 50,000 dollars per year? Yes or no?") as well as feature names and descriptions E_{name} , 200 E_{desc} (e.g., "hours-per-week": "the hours an individual has reported to work per week"). Similar to other works (Hegselmann et al., 2023), we serialize the sample data as in-context demonstration, 202 giving hints on the scale and format of the data. Given the data x, serialization is applied as: 203

Serialize(
$$\mathbf{x}, E_{\text{name}}$$
) = " e_{name}^1 is \mathbf{x}^1 . $\cdots e_{\text{name}}^d$ is \mathbf{x}^d .", (1)

where the superscript represents the vector's index value. 206

Operation instruction. This component guides the LLM on possible operations for feature discovery (see red part in Figure 2). It encourages the LLM to search only for feasibly-generated features, preventing erroneous behaviors (e.g., generating features that cannot be created from original data features or establishing ambiguous feature definitions). The operations used are as follows:

- Transformations: Transform the feature value with one of the following operators: absolute, logarithm, square root, sigmoid, or frequency.
- Numerical Operations: Conduct arithmetic operations from multiple numerical features.
 - Categorical Operations: Combine two categorical features to generate a new feature.

• Mixed-type Operations: Combine categorical and numerical features to generate a new one. For example, the model can discretize a numerical feature into a categorical one, allowing for categorical operations between the two features.

Diversity enforcement. Instead of concluding feature discovery with a single query, we aggregate
 features from multiple queries to find useful features. However, we also want to avoid the model from
 discovering duplicate features over multiple trials. To ensure diverse search, we provide additional
 instructions to prevent the LLM from selecting features identified in previous attempts (see teal part
 in Figure 2). We include descriptions via one-line pseudo-code, along with feature names, to prevent
 the LLM from simply renaming and selecting the same features. This component is integrated in all
 iterations except for the initial iteration.

Response instruction. This component includes instructions on how the LLM should format its response (see brown part in Figure 2). The format includes the type of operation, the feature's name, and a one-line pseudo-code necessary to regenerate the feature (e.g., "Numerical Operations | capital_diff | Subtract capital-loss from capital-gain to get the net capital difference"). Setting the response format facilitates easier parsing later on and also gives further evidence on each feature discovery, explaining why the particular feature was selected upon response.

The output text from the LLM prompt is parsed to generate the discovered feature. The generation process is automatically carried out using the LLM, which uses one-line pseudo-code and data to generate function code for producing the feature. The prompt used for automated generation is in the Appendix A.4.

237 238

216

217

218

219

3.2 Representation learning with discovered features

239 The discovered features are semantically related to the target downstream task based on LLM's prior 240 knowledge. We considered these features as ground-truth labels \hat{y} to define a pre-text task aligned 241 with the target downstream task. Although TST-LLM is agnostic to the choice of learning methods¹, 242 we adopt supervised contrastive learning for its generalizability to downstream tasks (Graf et al., 243 2021; Khosla et al., 2020). We define the projected representation of sample \mathbf{x}_i as $\mathbf{z}_i = g(f(\mathbf{x}_i))$, 244 where f is the encoder and q is the projection head. According to the literature (Khosla et al., 2020), 245 given a batched set of N_b samples with a pseudo label $\mathcal{B} = \{\mathbf{x}_i, \hat{y}_i\}_{i=1}^{N_b}$, the supervised contrastive 246 loss with a temperature au is defined as below. For numerical features among the discovered features, 247 we transformed them into discrete features using 1-dimensional k-means clustering with k = 10. 248

$$\mathcal{L}_{\text{SCL}} = -\frac{1}{|\mathcal{B}|} \sum_{i \in \mathcal{B}} \sum_{j \in \mathcal{B}, j \neq i} \mathbf{1}_{\hat{y}_i = \hat{y}_j} \log \frac{\exp(\mathbf{z}_i \cdot \mathbf{z}_j / \tau)}{\sum_{k \in \mathcal{B}} \mathbf{1}_{i \neq k} \exp(\mathbf{z}_i \cdot \mathbf{z}_k / \tau)}$$
(2)

(3)

In our framework, multiple features are discovered, and correspondingly, multiple labels are available for supervised contrastive learning. We utilize a multi-task learning approach to train the encoder, by defining a projection head for each ground-truth label and simultaneously training the models via supervised contrastive learning for each label. Specifically, given a set of M discovered features, $\hat{\mathcal{Y}} = {\hat{y}^1, \hat{y}^2, \dots, \hat{y}^M}$, we define a set of projection heads $\mathcal{G} = {g^1, g^2, \dots, g^M}$. Subsequently, the encoder f optimizes the following loss:

249 250 251

262

where each \mathcal{L}_{SCL}^m is a supervised contrastive loss computed with the respective projection head g^m for the corresponding label set \hat{y}^m .

 $\mathcal{L}_{\text{SCL-multi}} = \frac{1}{|M|} \sum_{m=1}^{M} \mathcal{L}_{\text{SCL}}^{m},$

Feature selection with minimum redundancy. Multi-task learning on all features generated by the LLM can be computationally heavy. Not all features are informative, and those that closely correlate with original features tend to limit their value as pre-texts. Furthermore, high correlation among the generated features could diminish the benefits of multi-task learning. To address this, we eliminate features that do not contribute meaningful information and carefully choose a diverse set of features with minimal redundancy, thereby reducing computation costs. (see Algorithm 1 in Appendix C.1).

¹See the Appendix E for the comparison with alternative learning methods (including reconstruction).

First, we define *uninformative* features as those with the lowest entropy values in the distribution. For example, if a feature is predominantly assigned to only one class across all samples (i.e., low entropy), the amount of information that can be learned from this feature is also limited. After calculating the entropy of each discovered feature, those with an entropy below a specific threshold (i.e., t_{ent}) are eliminated. The filtering threshold t_{ent} is set to 0.7, taking into account the entropy distribution of the entire feature set.

276 For the remaining features, the model selects a feature set that minimizes redundancy. The initial 277 choice is the feature with the smallest correlation to the original data. The remaining features are 278 then selected among the ones with the smallest correlation with the original data, including the 279 previously added features. This process is repeated until a predetermined number of features, M, 280 are selected. Cramer's V value (Cramér, 1999) is used to measure the correlation after discretizing all numerical features in the same manner as \hat{y} . This approach ensures that the selected features have 281 low correlation with the original data while maintaining diversity among the discovered features, 282 enabling efficient multi-task learning (see Table 1 for an analysis of feature diversity). Refer to the 283 Appendix G for example features generated and selected by our method. 284

285 286

287

292

4 EXPERIMENT

We evaluate TST-LLM across multiple tabular datasets with various downstream tasks. Through our experiments, we discuss which components of the model contributed to performance enhancements and how our model operates. Due to space constraints, full results, computational cost analysis, results with alternative learning objectives, and additional analyses can be found in the Appendix.

293 4.1 PERFORMANCE EVALUATION

Datasets. Our study used a total of 22 datasets to ensure a diverse range of downstream tasks in 295 terms of size and complexity including Adult (Asuncion & Newman, 2007), Balance-scale (Siegler, 296 1994), Bank (Moro et al., 2014), Blood (Yeh et al., 2009), Car (Kadra et al., 2021), Communities (Red-297 mond, 2009), Credit-g (Kadra et al., 2021), Diabetes (Smith et al., 1988), Eucalyptus (Bulloch et al., 298 1991), Forest-fires (Cortez & Morais, 2008), Heart (fedesoriano, 2021), Junglechess (van Rijn & Vis, 299 2014), Myocardial (Golovenkin et al., 2020), Tic-tac-toe (Aha, 1991), Vehicle (Mowforth & Shep-300 herd), Bike (Fanaee-T, 2013), Crab (Sidhu, 2021), Housing (Pace & Barry, 1997), Insurance (Datta, 301 2020), Wine (Cortez & Reis, 2009), Sequence-type, and Solution-mix. Descriptive statistics and task 302 descriptions for each dataset are available in the Appendix A.1 and B. Among them, 15 datasets were used for classification problems and 7 for regression problems; Two datasets-Sequence-type and 303 Solution-mix-are synthetic, ensuring they are not included in the LLM's pre-training corpus. 304

305 **Baselines.** Our model was compared to seven baselines, all of which were trained under the same 306 unsupervised setting as our experimental setup. (1) Raw Data: Uses the data as-is for the downstream 307 task without any representation learning; (2) AutoEncoder (Baldi, 2012): Utilizes a pre-text objective 308 that projects data into embeddings and reconstructs the original data; (3) SimSiam (Chen & He, 309 2021): Trains to minimize the embedding distance between a sample and its augmented version using 310 a siamese network structure; (4) SCARF (Bahri et al., 2022): Employs self-supervised contrastive learning to train augmentation-invariant embeddings. Augmentations involve corrupting some 311 columns of a sample by drawing from their marginal distributions; (5) STAB (Hajiramezanali et al., 312 2022): Similar to SimSiam but performs augmentation-free representation learning through stochastic 313 regularization; (6) STUNT (Nam et al., 2023b): Creates self-generated tasks based on clustering to 314 facilitate learning through meta-learning; (7) LFR (Sui et al., 2024): Iteratively learns the target of a 315 pre-text task and the encoder using a random data projector. Implementation details for all baselines 316 followed the original works, except that the encoder architecture was standardized. Detailed settings 317 can be found in the Appendix C.2. 318

Implementation details. TST-LLM currently employs GPT-3.5 as the LLM backbone for feature discovery but it can be combined with other LLMs. During LLM generation, the temperature was set to 0.5 and the top-p value was set to the API's default of 1. The discovery process generated five features per trial, with the number of trials set at 40. The number of serialized samples included in the prompt was set to a maximum of 20, as allowed by the prompt limit. The number of selected features M was set to 20. Effects of hyper-parameter M are discussed in the section 4.3 and an



338 Figure 3: Win matrices comparing self-supervised tabular learning methods against each other with 339 (a) linear model and (b) non-parametric classifier. Self-supervised tabular learning methods are aligned on the x-axis and the y-axis while the numbers represent the winning ratio of the x-axis model 340 against the y-axis model. Full results are reported in the Appendix H. 341

analysis of the number of trials is presented in the Appendix F. The structure of the encoder was consistent with the baselines, configured as a 2-layer MLP with 1024 dimensions, and the projection head consisted of a single linear layer. Training utilized the Adam optimizer with a learning rate of 1e-4, a batch size of 128, and 1000 training iterations. For information on computing resources and computational complexity, refer to the Appendix D.

349 **Evaluation.** After training, we fixed the learned embeddings, and the evaluation is performed with two downstream task classifiers: (1) Linear model: This can be either logistic regression for 350 classification tasks or linear regression for regression tasks. This method assesses how linearly 351 separable the classes are in the embedding; (2) Non-parametric classifier: This involves fitting a 352 weighted k-NN module to the downstream classification task. We run evaluations with two different 353 settings: k = 3 and 5. This method evaluates how well the embeddings form coherent local clusters. 354 For performance metrics, AUROC is used for classification tasks (one-versus-all for multi-class 355 settings), and RMSE for regression tasks. Experiments were run with 3 different random seeds, and 356 the average values were reported. 357

To facilitate straightforward comparison across datasets, we adopted a win matrix from existing 358 literature (Bahri et al., 2022). The win matrix calculates the ratio over the number of times each 359 method *i* outperforms another method *j* across the datasets, excluding ties: 360

$$W[i, j] = \frac{\sum_{k \in \text{Datasets}} \mathbb{I}[\text{Performance}(i, k) > \text{Performance}(j, k)]}{\sum_{k \in \text{Datasets}} \mathbb{I}[\text{Performance}(i, k) \neq \text{Performance}(j, k)]},$$
(4)

where Performance(i, k) denotes the performance of method *i* on dataset *k*.

365 **Results.** Figure 3 compares the performance of self-supervised baselines and TST-LLM against 366 each other using win matrices. For all baselines, the average win ratio is 84% for the linear model and 65% for the non-parametric classifier, demonstrating TST-LLM's superiority; This gives a 368 strong evidence that task-specific pre-text tasks lead to the latent representations that readily form 369 decision boundaries for the target downstream task in the case of the linear model. At the same time, 370 evaluations with a non-parametric classifier indicate that our pre-text tasks effectively extract and utilize information from existing features to enable clustering.

373 4.2 ABLATION STUDY

342 343

344

345

346

347

348

361 362

364

367

371

372

374

375 We conducted an ablation study to evaluate the contribution of each component in our model. We assessed two primary components: discovering features from the downstream task's description and 376 training the encoder with multi-task contrastive learning. We defined the following ablations by 377 removing or modifying each component: (1) Top-1 selection: Only the top-1 feature, which has



Figure 4: Win matrices comparing our full model and its ablations against each other with (a) linear model and (b) non-parametric classifier. The numbers represent the winning ratio.

393 the least redundancy with the original data among the discovered features, is used; (2) Random-1 394 selection: Same as Top-1 selection, a single head is used for training, yet the label used for supervised 395 contrastive learning is randomly changed to one of the discovered features in each iteration; (3) 396 Random feature discovery: Instead of using the LLM for feature discovery, we expand features 397 using operations commonly employed in traditional feature engineering work (Zhang et al., 2023), and then randomly select M features. Representation learning is subsequently conducted with 398 these selected features, identical to our original model's approach; (4) Without learning: Instead 399 of performing representation learning, features discovered through feature discovery are directly 400 concatenated with the original data and used as is; (5) Without feature selection: All discovered 401 features are used in representation learning without undergoing the feature selection process. 402

403 Figure 4 shows the degree of performance degradation in each ablation study. We find that every ablation led to a negative effect on performance, underscoring the contribution of all tested compo-404 nents. Specifically, using multiple features for multi-task learning, rather than relying on a single 405 feature (i.e., Top-1 selection) or alternating features for single-task learning (i.e., Random-1 selection), 406 provided an ensemble effect that enhanced performance. Even without training, merely concatenating 407 features that are relevant to the actual label facilitated the formation of effective local clusters with the 408 non-parametric classifier. By conducting training with a pre-text task, TST-LLM could further obtain 409 embeddings that are linearly separable among the labels. In addition, selecting features does not 410 significantly differ in performance from using all features without feature selection, which suggests 411 that our selection strategy leads to efficient learning (see Table 1 for comparison on computational 412 complexity of using all features).

413 414

415

390

391 392

4.3 ANALYSIS & DISCUSSION

416 How informative are the discovered features for the downstream task? When training 417 TST-LLM, we utilize the features that have been identified through the LLM. To see how well these 418 pre-text tasks align with actual downstream tasks, we computed the average increase ratio of mutual information between the discovered features and downstream task's labels compared to the original 419 features. The ratio is computed as a percentage for each dataset. According to Figure 5a, for most 420 datasets, the discovered features show a stronger correlation with the labels than the original data. 421 This suggests that our pre-text tasks are more closely aligned with the actual downstream tasks, 422 than the models trained solely on the original data. We also observed that datasets with a higher 423 increase ratio also demonstrated a greater performance improvement in our model compared to the 424 Raw Data model (Spearman correlation 0.52). When evaluating our model compared to the raw 425 data model over datasets with positive and negative increase ratios, the positive set showed a 25.8% 426 greater improvement in performance (average 8.1% increase in the positive set vs. 5.6% increase 427 in the negative set). Although some datasets showed a decrease in average mutual information, most 428 of them exhibited a high standard deviation in mutual information (see Appendix H.5 for full results), 429 where multi-task learning using a variety of features could be helpful.

- 430
- **How diverse are the features used for our pre-text task?** We applied two strategies to ensure the diversity of discovered features coming from the LLM and pre-text tasks for representation



Figure 5: (a) Average increase ratio of mutual information between the discovered features and ground-truth labels compared to the original features. The ratio is reported as a percentage for each dataset; (b) Hyper-parameter analysis on the number of selected features M. Average decrease ratios from our model's settings (i.e., M = 20) across both classification and regression tasks are reported.

learning. One strategy involved adding a diversity enforcement component within the LLM's prompt to avoid selecting the previously selected features, and the other aimed to minimize redundancy 452 among selected features. To verify the effectiveness of these methodologies on feature diversity, we 453 conducted additional experiments. We defined three ablation scenarios: (1) No diversity enforcement 454 & No selection strategy: using features without applying the two strategies; (2) No selection strategy: 455 using the diversity enforcement component but not conducting feature selection; (3) With entropy-456 based filtering: applying only entropy-based filtering as the selection strategy.

457 We compared each ablation using three evaluation metrics. The first metric is Cramer's V 458 value (Cramér, 1999) between features, where a higher value indicates a greater number of highly 459 correlated features, implying lower diversity. The second metric is the percentage change in perfor-460 mance across all datasets compared to the proposed full model. The final metric is the time cost ratio 461 for running the model. According to the results in Table 1, models with lower diversity are inefficient 462 both in terms of performance and time cost. 463

Table 1: Ablation study results on feature diversity. Feature diversity is evaluated using the average 464 Cramer's V value across features, with the standard deviation noted. Performance change is computed 465 as an averaged change ratio in percentage across all datasets compared to the proposed full model. 466

Ablation for feature diversity	Cramer's V	Performance change (%)	Time cost ratio
No diversity enforcement & No selection strategy	0.24±0.11	-0.17±0.24	5.62
No selection strategy	$0.13 {\pm} 0.06$	-0.50 ± 0.38	4.84
With entropy-based filtering	$0.09{\pm}0.05$	-0.05 ± 0.15	2.78
Full model	$0.07 {\pm} 0.03$	$0.00{\pm}0.00$	1.00

473

445

446

447

448 449 450

451

474 **Does hyper-parameter** M **affect the performance?** TST-LLM has a hyper-parameter, M, which 475 represents the number of features discovered for the pre-text task. To investigate the impact of M on performance, we conducted experiments using M = 10, 20, 30, and all features (i.e., M = all) for 476 the pretext task. The results, presented in Figure 5b, include the average decrease ratio in performance 477 from our model's settings across all classification and regression datasets. The performance of 478 TST-LLM is insensitive to M when M is set bigger than 10, allowing for flexibility in choosing the 479 number of features to optimize computational efficiency. Based on our findings, we selected M = 20, 480 which delivered the best performance without imposing a computational burden. 481

482

5 CONCLUSION

483 484

We introduced TST-LLM, a representation learning method that creates pre-text tasks that are tailored 485 to downstream task objectives using an LLM. TST-LLM leverages the prior knowledge and reasoning abilities of the LLM to determine how to combine original data features into informative features
 based on natural language descriptions of downstream tasks and feature descriptions. The combined
 features, after undergoing a feature selection process to minimize redundancy, serve as ground truth labels for the pre-text tasks in representation learning. Extensive analysis confirms that our
 methodology can identify diverse and task-aligned features, and as a result consistently achieves
 outstanding performance across various downstream tasks.

492 Future work and broader impact. Our method relies on LLM for feature discovery, which may 493 not yield optimal results for tasks that the LLM is unfamiliar with. To mitigate this, one could 494 consider optimizing alongside traditional self-supervised representation learning objectives in the 495 tabular domain, such as reconstruction or contrastive learning. Alternatively, one could consider 496 calibrating the discovered features through human feedback. In terms of the impact, TST-LLM 497 facilitates easy learning through task-aligned pre-text tasks with the desired downstream task 498 objective, when these goals can be articulated through text. This adaptability renders it suitable for 499 a variety of real-world scenarios, such as in the healthcare and financial sectors. We believe this 500 work provides a new perspective on the integration of LLMs into the tabular learning domain.

References

501 502

503

504

505

524

525

526

530

531

532

533

- David Aha. Tic-Tac-Toe Endgame. UCI Machine Learning Repository, 1991. DOI: https://doi.org/10.24432/C5688J.
- Rohan Anil, Andrew M Dai, Orhan Firat, Melvin Johnson, Dmitry Lepikhin, Alexandre Passos, Siamak Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng Chen, et al. Palm 2 technical report. *arXiv preprint arXiv:2305.10403*, 2023.

Mahmoud Assran, Mathilde Caron, Ishan Misra, Piotr Bojanowski, Florian Bordes, Pascal Vincent,
 Armand Joulin, Mike Rabbat, and Nicolas Ballas. Masked siamese networks for label-efficient
 learning. In *European Conference on Computer Vision*, pp. 456–473. Springer, 2022.

- Arthur Asuncion and David Newman. Uci machine learning repository, 2007.
- 515 Dara Bahri, Heinrich Jiang, Yi Tay, and Donald Metzler. Scarf: Self-supervised contrastive learning 516 using random feature corruption. In *International Conference on Learning Representations*, 2022.
- Pierre Baldi. Autoencoders, unsupervised learning, and deep architectures. In *Proceedings of ICML workshop on unsupervised and transfer learning*, pp. 37–49. JMLR Workshop and Conference Proceedings, 2012.
- Randall Balestriero, Mark Ibrahim, Vlad Sobal, Ari Morcos, Shashank Shekhar, Tom Goldstein,
 Florian Bordes, Adrien Bardes, Gregoire Mialon, Yuandong Tian, et al. A cookbook of self supervised learning. *arXiv preprint arXiv:2304.12210*, 2023.
 - BT Bulloch et al. Eucalyptus species selection for soil conservation in seasonally dry hill countrytwelfth year assessment. *New Zealand journal of forestry science*, 21(1):10–31, 1991.
- Mathilde Caron, Ishan Misra, Julien Mairal, Priya Goyal, Piotr Bojanowski, and Armand Joulin.
 Unsupervised learning of visual features by contrasting cluster assignments. *Advances in neural information processing systems*, 33:9912–9924, 2020.
 - Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pp. 1597–1607. PMLR, 2020.
- Xinlei Chen and Kaiming He. Exploring simple siamese representation learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 15750–15758, 2021.
- Cerdeira A.-Almeida F. Matos T. Cortez, Paulo and J. Reis. Wine Quality. UCI Machine Learning Repository, 2009. DOI: https://doi.org/10.24432/C56S3T.
- 539 Paulo Cortez and Anbal Morais. Forest Fires. UCI Machine Learning Repository, 2008. DOI: https://doi.org/10.24432/C5D88D.

540 541	Harald Cramér. Mathematical methods of statistics, volume 26. Princeton university press, 1999.		
542	Anirban Datta US Health Insurance Dataset Kaggle 2020 kag-		
543	gle.com/datasets/teertha/ushealthinsurancedataset.		
544			
545	Tuan Dinh, Yuchen Zeng, Ruisu Zhang, Ziqian Lin, Michael Gira, Shashank Rajput, Jy-yong		
546	Sohn, Dimitris Papailiopoulos, and Kangwook Lee. Lift: Language-interfaced fine-tuning for		
547	non-language machine learning tasks. Advances in Neural Information Processing Systems, 35:		
5/18	11/63–11/84, 2022.		
540	Hadi Fanaee-T. Bike Sharing. UCI Machine Learning Repository. 2013. DOI:		
550	https://doi.org/10.24432/C5W894.		
551	fedesoriano. Heart Failure Prediction Dataset. Kaggle, 2021. kaggle.com/fedesoriano/heart-failure-		
552 553	prediction.		
554	Tianyu Gao, Xingcheng Yao, and Danqi Chen. Simcse: Simple contrastive learning of sentence		
555 556	embeddings. In Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, pp. 6894–6910, 2021.		
557	0/11 /		
558	Zhabiz Gharibshah and Xingquan Zhu. Local contrastive feature learning for tabular data. In		
550	Proceedings of the 31st ACM International Conference on Information & Knowledge Management,		
559	pp. 3963–3967, 2022.		
561	Spyros Gidaris Praveer Singh and Nikos Komodakis. Unsupervised representation learning by		
562	predicting image rotations. In International Conference on Learning Representations, 2018.		
562			
564	S.E. Golovenkin, V.A. Shulman, D.A. Rossiev, P.A. Shesternya, S.Yu. Nikulina, Yu.V. Orlova, and		
565	V.F. Voino-Yasenetsky. Myocardial infarction complications. UCI Machine Learning Repository,		
566	2020. DOI: https://doi.org/10.24432/C53P5M.		
567	Florian Graf, Christoph Hofer, Marc Niethammer, and Roland Kwitt. Dissecting supervised con-		
568	trastive learning. In International Conference on Machine Learning, pp. 3821–3830. PMLR,		
569	2021.		
570			
571	Jean-Bastien Grill, Florian Strub, Florent Altche, Corentin Tallec, Pierre Richemond, Elena		
572	buchatskaya, Cari Doelsch, Bernardo Avna Pires, Zhaonan Guo, Monaininad Gheshiagin Azar, et al. Bootstrap your own latent a new approach to self supervised learning. Advances in neural		
573	information processing systems 33.21271–21284 2020		
574	<i>uyormanon processing systems, 55.21271 21201, 2020.</i>		
575	Ehsan Hajiramezanali, Nathaniel Lee Diamant, Gabriele Scalia, and Max W Shen. Stab: Self-		
576	supervised learning for tabular data. In NeurIPS 2022 First Table Representation Work		
577	2022.		
578	Sungwon Han Sungwon Park Sungkyu Park Sundong Kim and Meavoung Cha. Mitigating		
579	embedding and class assignment mismatch in unsupervised image classification. In <i>Furgage</i>		
580	Conference on Computer Vision, pp. 768–784. Springer, 2020.		
581			
582	Sungwon Han, Jinsung Yoon, Sercan O Arik, and Tomas Pfister. Large language models can		
583	automatically engineer features for few-shot tabular learning. arXiv preprint arXiv:2404.09491,		
584	2024.		
585	Stefan Hegselmann, Alejandro Buendia, Hunter Lang, Monica Agrawal, Xiaovi Jiang, and David		
586	Sontag. Tablim: Few-shot classification of tabular data with large language models. In <i>International</i>		
587	Conference on Artificial Intelligence and Statistics, pp. 5549–5581. PMLR, 2023.		
588			
589	Noah Hollmann, Samuel Müller, Katharina Eggensperger, and Frank Hutter. Tabpfn: A transformer		
590	unat solves small tabular classification problems in a second. In <i>The Eleventh International</i>		
591	Conjerence on Learning Representations, 2025.		
592	Edward J Hu, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen,		
593	et al. Lora: Low-rank adaptation of large language models. In International Conference on Learning Representations, 2021.		

- Arlind Kadra, Marius Lindauer, Frank Hutter, and Josif Grabocka. Well-tuned simple nets excel on tabular datasets. *Advances in neural information processing systems*, 34:23928–23941, 2021.
- Jacob Devlin Ming-Wei Chang Kenton and Lee Kristina Toutanova. Bert: Pre-training of deep
 bidirectional transformers for language understanding. In *Proceedings of NAACL-HLT*, pp. 4171–
 4186, 2019.
- Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron Maschinot, Ce Liu, and Dilip Krishnan. Supervised contrastive learning. *Advances in neural information processing systems*, 33:18661–18673, 2020.
- Dahun Kim, Donghyeon Cho, Donggeun Yoo, and In So Kweon. Learning image representations
 by completing damaged jigsaw puzzles. In 2018 IEEE Winter Conference on Applications of
 Computer Vision (WACV), pp. 793–802. IEEE, 2018.
- Haokun Liu, Derek Tam, Mohammed Muqeeth, Jay Mohta, Tenghao Huang, Mohit Bansal, and
 Colin A Raffel. Few-shot parameter-efficient fine-tuning is better and cheaper than in-context
 learning. Advances in Neural Information Processing Systems, 35:1950–1965, 2022.
- Himangi Mittal, Pedro Morgado, Unnat Jain, and Abhinav Gupta. Learning state-aware visual representations from audible interactions. *Advances in Neural Information Processing Systems*, 35: 23765–23779, 2022.
- Sérgio Moro, Paulo Cortez, and Paulo Rita. A data-driven approach to predict the success of bank
 telemarketing. *Decision Support Systems*, 62:22–31, 2014.
- 616
 617 Pete Mowforth and Barry Shepherd. Statlog (Vehicle Silhouettes). UCI Machine Learning Repository. DOI: https://doi.org/10.24432/C5HG6N.
- Jaehyun Nam, Woomin Song, Seong Hyeon Park, Jihoon Tack, Sukmin Yun, Jaehyung Kim, and
 Jinwoo Shin. Semi-supervised tabular classification via in-context learning of large language
 models. In *Workshop on Efficient Systems for Foundation Models (CML2023)*, 2023a.
- Jaehyun Nam, Jihoon Tack, Kyungmin Lee, Hankook Lee, and Jinwoo Shin. Stunt: Few-shot tabular
 learning with self-generated tasks from unlabeled tables. In *The Eleventh International Conference on Learning Representations*, 2023b.
- Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018.
- 629 OpenAI. Gpt-4 technical report, 2023.

610

622

- Andrew Owens and Alexei A Efros. Audio-visual scene analysis with self-supervised multisensory features. In *Proceedings of the European conference on computer vision (ECCV)*, pp. 631–648, 2018.
- R Kelley Pace and Ronald Barry. Sparse spatial autoregressions. *Statistics & Probability Letters*, 33 (3):291–297, 1997.
- Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, Jeffrey Dean, and Sanjay Ghemawat. Language models are unsupervised multitask learners. In OSDI'04: Sixth Symposium on Operating System Design and Implementation, pp. 137–150, 2019.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual
 models from natural language supervision. In *International conference on machine learning*, pp.
 8748–8763. PMLR, 2021.
- Michael Redmond. Communities and Crime. UCI Machine Learning Repository, 2009. DOI: https://doi.org/10.24432/C53W3X.
- 647 Gursewak Singh Sidhu. Crab age prediction, 2021. URL https://www.kaggle.com/dsv/ 2834512.

648 649 650	R. Siegler. Balance Scale. UCI Machine Learning Repository, 1994. DOI: https://doi.org/10.24432/C5488X.
651 652 653 654	Jack W Smith, James E Everhart, WC Dickson, William C Knowler, and Robert Scott Johannes. Using the adap learning algorithm to forecast the onset of diabetes mellitus. In <i>Proceedings of the annual symposium on computer application in medical care</i> , pp. 261. American Medical Informatics Association, 1988.
655 656 657	Gowthami Somepalli, Avi Schwarzschild, Micah Goldblum, C Bayan Bruss, and Tom Goldstein. Saint: Improved neural networks for tabular data via row attention and contrastive pre-training. In <i>NeurIPS 2022 First Table Representation Workshop</i> , 2022.
658 659 660 661	Yi Sui, Tongzi Wu, Jesse C Cresswell, Ga Wu, George Stein, Xiao Shi Huang, Xiaochen Zhang, and Maksims Volkovs. Self-supervised representation learning from random data projectors. In <i>The Twelfth International Conference on Learning Representations</i> , 2024.
662 663 664	Michael Tschannen, Josip Djolonga, Paul K Rubenstein, Sylvain Gelly, and Mario Lucic. On mutual information maximization for representation learning. In <i>International Conference on Learning Representations</i> , 2019.
665 666 667	Talip Ucar, Ehsan Hajiramezanali, and Lindsay Edwards. Subtab: Subsetting features of tabular data for self-supervised representation learning. <i>Advances in Neural Information Processing Systems</i> , 34:18853–18865, 2021.
668 669 670	Jan N van Rijn and Jonathan K Vis. Endgame analysis of dou shou qi. <i>ICGA Journal</i> , 37(2):120–124, 2014.
671 672	Zifeng Wang and Jimeng Sun. Transtab: Learning transferable tabular transformers across tables. <i>Advances in Neural Information Processing Systems</i> , 35:2902–2915, 2022.
674 675	Zifeng Wang, Chufan Gao, Cao Xiao, and Jimeng Sun. Anypredict: Foundation model for tabular prediction. <i>arXiv preprint arXiv:2305.12081</i> , 2023.
676 677 678	Xin Wen, Bingchen Zhao, Anlin Zheng, Xiangyu Zhang, and Xiaojuan Qi. Self-supervised visual representation learning with semantic grouping. <i>Advances in neural information processing systems</i> , 35:16423–16438, 2022.
679 680 681 682 683	Jing Wu, Suiyao Chen, Qi Zhao, Renat Sergazinov, Chen Li, Shengjie Liu, Chongchao Zhao, Tianpei Xie, Hanqing Guo, Cheng Ji, et al. Switchtab: Switched autoencoders are effective tabular learners. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 38, pp. 15924–15933, 2024.
684 685 686	Zhirong Wu, Yuanjun Xiong, Stella X Yu, and Dahua Lin. Unsupervised feature learning via non- parametric instance discrimination. In <i>Proceedings of the IEEE conference on computer vision</i> <i>and pattern recognition</i> , pp. 3733–3742, 2018.
687 688	I-Cheng Yeh, King-Jang Yang, and Tao-Ming Ting. Knowledge discovery on rfm model using bernoulli sequence. <i>Expert Systems with applications</i> , 36(3):5866–5871, 2009.
690 691 692	Jinsung Yoon, Yao Zhang, James Jordon, and Mihaela van der Schaar. Vime: Extending the success of self-and semi-supervised learning to tabular domain. <i>Advances in Neural Information Processing Systems</i> , 33:11033–11043, 2020.
693 694 695	Jure Zbontar, Li Jing, Ishan Misra, Yann LeCun, and Stéphane Deny. Barlow twins: Self-supervised learning via redundancy reduction. In <i>International conference on machine learning</i> , pp. 12310–12320. PMLR, 2021.
696 697 698 699	Richard Zhang, Phillip Isola, and Alexei A Efros. Split-brain autoencoders: Unsupervised learning by cross-channel prediction. In <i>Proceedings of the IEEE conference on computer vision and pattern recognition</i> , pp. 1058–1067, 2017.
700 701	Tianping Zhang, Zheyu Aqa Zhang, Zhiyuan Fan, Haoyan Luo, Fengyuan Liu, Qian Liu, Wei Cao, and Li Jian. Openfe: Automated feature generation with expert-level performance. In <i>International Conference on Machine Learning</i> , pp. 41880–41901. PMLR, 2023.

702	Bingzhao Zhu, Xingjian Shi, Nick Erickson, Mu Li, George Karypis, and Mahsa Shoaran.	Xtab:
703	Cross-table pretraining for tabular transformers. arXiv preprint arXiv:2305.06090, 2023.	
704		
705		
706		
707		
708		
709		
710		
711		
712		
713		
714		
715		
716		
717		
718		
719		
720		
721		
722		
723		
724		
725		
726		
727		
728		
729		
730		
731		
732		
733		
734		
730		
730		
729		
730		
7/0		
7/1		
742		
7/3		
744		
745		
746		
747		
748		
749		
750		
751		
752		
753		
754		
755		

756 757

758

759

APPENDIX

A FULL PROMPT EXAMPLES

760 761 A.1 TASK DESCRIPTION FOR EACH DATASET

This section presents the downstream task descriptions corresponding to the dataset used for evaluation. TST-LLM uses these text descriptions to perform task-relevant feature discovery. Each
description is defined by referencing the dataset's original source or previous works (Hegselmann
et al., 2023; Han et al., 2024). For classification tasks, answer class candidates were provided.

766 767

Table 2: Downstream task's description of each dataset used for feature discovery.

Data	Downstream task's description
Adult	Does this person earn more than 50000 dollars per year? Yes or no?
Balance-scale	Which direction does the balance scale tip to? Right, left, or balanced?
Bank	Does this client subscribe to a term deposit? Yes or no?
Blood	Did the person donate blood? Yes or no?
Car	How would you rate the decision to buy this car?
a	Unacceptable, acceptable, good or very good?
Communities	How high will the rate of violent crimes per 100K population be in this area.
Cradit a	Low, medium, or mign?
Jieun-g Diabetes	Does this period fective a credit? Tes of no?
Fucalyptus	How good is this Eucalyntus species for soil conservation
Eucuryptus	in the specified location? None, low average, good, or best?
Forest-fires	Estimate the burned area of forest fires from given information.
Heart	Does the coronary angiography of this patient show a heart disease? Yes or no?
Junglechess	Which player wins this two pieces endgame of Jungle Chess? Black, white or draw?
Myocardial	Does the myocardial infarction complications data of this patient show
	chronic heart failure? Yes or no?
Tic-tac-toe	Will the first player (player x) win the game? Positive or negative?
Vehicle	What kind of vehicle is the given silhouette information about? Bus, opel, saab, or van?
Bike	Estimate the count of total rental bikes from given information.
Crab	Estimate the age of the crab from given information.
Housing	Estimate the nouse price from given information.
Wine	Estimate the matvioual medical cost of this patient officer by health insurance. Estimate the wine quality on a scale from 0 to 10 from given information
Sequence-type	What is the type of following sequence? Arithmetic geometric fibonacci or collatz?
Solution-mix	Given the volumes and concentrations of four solutions.
bolution min	calculate the percent concentration of the mixed solution after mixing them.

810 A.2 FULL PROMPT EXAMPLE FOR FEATURE DISCOVERY 811

The following is an example of a prompt used for feature discovery on the Adult dataset. For the initial query in the LLM, a prompt without a diversity enforcement component, as shown in Figure 6, was used as there is no information from previous iterations. For subsequent iterations, a prompt with a diversity enforcement component in Figure 7 was used.

a engineer. Given the task description and the list of features and data examples, you are w column for the data which is informative to solve the task. his person earn more than 50000 dollars per year? Yes or no? of an individual (numerical variable within range [17, 90]) rry: country of origin for an individual (categorical variable with categories [United-States, foland-Netherlands]) workclass is Private. fnlwgt is 123807. education is HS-grad. educational-num is 9. is Separated. occupation is Adm-clerical. relationship is Unmarried. race is Black. gender upital-gain is 0. capital-loss is 0. hours-per-week is 40. native-country is United-States. vorkclass is Private. fnlwgt is 208302. education is HS-grad. educational-num is 9. is Married-civ-spouse. occupation is Sales. relationship is Husband. race is White. le. capital-gain is 0. capital-loss is 0. hours-per-week is 36. native-country is United-States.
 w column for the data which is informative to solve the task. his person earn more than 50000 dollars per year? Yes or no? e of an individual (numerical variable within range [17, 90]) try: country of origin for an individual (categorical variable with categories [United-States, toland-Netherlands]) vorkclass is Private. fnlwgt is 123807. education is HS-grad. educational-num is 9. s is Separated. occupation is Adm-clerical. relationship is Unmarried. race is Black. gender upital-gain is 0. capital-loss is 0. hours-per-week is 40. native-country is United-States. vorkclass is Private. fnlwgt is 208302. education is HS-grad. educational-num is 9. s is Married-civ-spouse. occupation is Sales. relationship is Husband. race is White. le. capital-gain is 0. capital-loss is 0. hours-per-week is 36. native-country is United-States.
his person earn more than 50000 dollars per year? Yes or no? e of an individual (numerical variable within range [17, 90]) rry: country of origin for an individual (categorical variable with categories [United-States, oland-Netherlands]) vorkclass is Private. fnlwgt is 123807. education is HS-grad. educational-num is 9. s is Separated. occupation is Adm-clerical. relationship is Unmarried. race is Black. gender upital-gain is 0. capital-loss is 0. hours-per-week is 40. native-country is United-States. vorkclass is Private. fnlwgt is 208302. education is HS-grad. educational-num is 9. s is Married-civ-spouse. occupation is Sales. relationship is Husband. race is White. le. capital-gain is 0. capital-loss is 0. hours-per-week is 36. native-country is United-States.
his person earn more than 50000 dollars per year? Yes or no? e of an individual (numerical variable within range [17, 90]) rry: country of origin for an individual (categorical variable with categories [United-States, oland-Netherlands]) vorkclass is Private. fnlwgt is 123807. education is HS-grad. educational-num is 9. s is Separated. occupation is Adm-clerical. relationship is Unmarried. race is Black. gender upital-gain is 0. capital-loss is 0. hours-per-week is 40. native-country is United-States. vorkclass is Private. fnlwgt is 208302. education is HS-grad. educational-num is 9. s is Married-civ-spouse. occupation is Sales. relationship is Husband. race is White. le. capital-gain is 0. capital-loss is 0. hours-per-week is 36. native-country is United-States.
e of an individual (numerical variable within range [17, 90]) try: country of origin for an individual (categorical variable with categories [United-States, oland-Netherlands]) vorkclass is Private. fnlwgt is 123807. education is HS-grad. educational-num is 9. s is Separated. occupation is Adm-clerical. relationship is Unmarried. race is Black. gender upital-gain is 0. capital-loss is 0. hours-per-week is 40. native-country is United-States. vorkclass is Private. fnlwgt is 208302. education is HS-grad. educational-num is 9. s is Married-civ-spouse. occupation is Sales. relationship is Husband. race is White. le. capital-gain is 0. capital-loss is 0. hours-per-week is 36. native-country is United-States.
 of an individual (numerical variable within range [17, 90]) try: country of origin for an individual (categorical variable with categories [United-States, ioland-Netherlands]) vorkclass is Private. fnlwgt is 123807. education is HS-grad. educational-num is 9. s is Separated. occupation is Adm-clerical. relationship is Unmarried. race is Black. gender upital-gain is 0. capital-loss is 0. hours-per-week is 40. native-country is United-States. vorkclass is Private. fnlwgt is 208302. education is HS-grad. educational-num is 9. s is Married-civ-spouse. occupation is Sales. relationship is Husband. race is White. le. capital-gain is 0. capital-loss is 0. hours-per-week is 36. native-country is United-States.
try: country of origin for an individual (categorical variable with categories [United-States, foland-Netherlands]) workclass is Private. fnlwgt is 123807. education is HS-grad. educational-num is 9. is Separated. occupation is Adm-clerical. relationship is Unmarried. race is Black. gender upital-gain is 0. capital-loss is 0. hours-per-week is 40. native-country is United-States. workclass is Private. fnlwgt is 208302. education is HS-grad. educational-num is 9. s is Married-civ-spouse. occupation is Sales. relationship is Husband. race is White. le. capital-gain is 0. capital-loss is 0. hours-per-week is 36. native-country is United-States.
vorkclass is Private. fnlwgt is 123807. education is HS-grad. educational-num is 9. s is Separated. occupation is Adm-clerical. relationship is Unmarried. race is Black. gender upital-gain is 0. capital-loss is 0. hours-per-week is 40. native-country is United-States. vorkclass is Private. fnlwgt is 208302. education is HS-grad. educational-num is 9. s is Married-civ-spouse. occupation is Sales. relationship is Husband. race is White. le. capital-gain is 0. capital-loss is 0. hours-per-week is 36. native-country is United-States.
vorkclass is Private. fnlwgt is 123807. education is HS-grad. educational-num is 9. is Separated. occupation is Adm-clerical. relationship is Unmarried. race is Black. gender upital-gain is 0. capital-loss is 0. hours-per-week is 40. native-country is United-States. vorkclass is Private. fnlwgt is 208302. education is HS-grad. educational-num is 9. s is Married-civ-spouse. occupation is Sales. relationship is Husband. race is White. le. capital-gain is 0. capital-loss is 0. hours-per-week is 36. native-country is United-States.
vorkclass is Private. fnlwgt is 123807. education is HS-grad. educational-num is 9. s is Separated. occupation is Adm-clerical. relationship is Unmarried. race is Black. gender upital-gain is 0. capital-loss is 0. hours-per-week is 40. native-country is United-States. vorkclass is Private. fnlwgt is 208302. education is HS-grad. educational-num is 9. s is Married-civ-spouse. occupation is Sales. relationship is Husband. race is White. le. capital-gain is 0. capital-loss is 0. hours-per-week is 36. native-country is United-States.
vorkclass is Private. fnlwgt is 123807. education is HS-grad. educational-num is 9. s is Separated. occupation is Adm-clerical. relationship is Unmarried. race is Black, gender ipital-gain is 0. capital-loss is 0. hours-per-week is 40. native-country is United-States. vorkclass is Private. fnlwgt is 208302. education is HS-grad. educational-num is 9. s is Married-civ-spouse. occupation is Sales. relationship is Husband. race is White. le. capital-gain is 0. capital-loss is 0. hours-per-week is 36. native-country is United-States.
vorkclass is Private. fnlwgt is 208302. education is HS-grad. educational-num is 9. s is Married-civ-spouse. occupation is Sales. relationship is Husband. race is White. le. capital-gain is 0. capital-loss is 0. hours-per-week is 36. native-country is United-States.
vorkclass is Private. fnlwgt is 208302. education is HS-grad. educational-num is 9. s is Married-civ-spouse. occupation is Sales. relationship is Husband. race is White. le. capital-gain is 0. capital-loss is 0. hours-per-week is 36. native-country is United-States.
vorkclass is Private. fnlwgt is 208302. education is HS-grad. educational-num is 9. s is Married-civ-spouse. occupation is Sales. relationship is Husband. race is White. le. capital-gain is 0. capital-loss is 0. hours-per-week is 36. native-country is United-States.
s is Married-civ-spouse. occupation is Sales. relationship is Husband. race is White. le. capital-gain is 0. capital-loss is 0. hours-per-week is 36. native-country is United-States.
le. capital-gain is 0. capital-loss is 0. hours-per-week is 36. native-country is United-States.
a of anorations below, concrete 5 new columns which are the most informative to
k using operations. Refer to the examples when generating features. Only use features
eature description. Note that multiple operations can be nested to generate a new column.
type of operations is as follows:
tions: Numerical features only. Transform the feature value with one of the following
arithm square root sigmoid or frequency (i.e. frequency of feature in the data)
Distribution of the second sec
Operations: Combine categorical feature and numerical feature to generate a new one.
Operations: Combine two categorical features to generate a new feature. For example, you
ondition to make a binary feature, indicating whether it follows the condition.
manaa far 5 navy aalumnay
sponse for 5 new columns:
Any reasons based on examples above why the following new feature would be helpful for
,
1: [Type of operation] New_column_name One line detailed pseudo code for generating
5.
5:
5:
5:
5:

Figure 6: Full prompt example for feature discovery in the Adult dataset (initial query without diversity enforcement).

- 859 860
- 861
- 862
- 863

864 You are a data engineer. Given the task description and the list of features and data examples, you are 865 making a new column for the data which is informative to solve the task. 866 867 Task: Does this person earn more than 50000 dollars per year? Yes or no? 868 Features: - age: the age of an individual (numerical variable within range [17, 90]) 870 ... 871 Examples: 872 age is 49. workclass is Private. fnlwgt is 123807. education is HS-grad. educational-num is 9. 873 marital-status is Separated. occupation is Adm-clerical. relationship is Unmarried. race is Black. gender 874 is Female. capital-gain is 0. capital-loss is 0. hours-per-week is 40. native-country is United-States. 875 ••• 876 Given a type of operations below, generate 5 new columns which are the most informative to 877 solve the task using operations. Refer to the examples when generating features. Only use features 878 listed in the feature description. Note that multiple operations can be nested to generate a new column. 879 The possible type of operations is as follows: - Transformations: Numerical features only. Transform the feature value with one of the following 881 operators: 882 absolute, logarithm, square root, sigmoid, or frequency (i.e., frequency of feature in the data). 883 - Numerical Operations: Numerical features only. Conduct arithmetic operation from multiple columns. 884 - Mixed-type Operations: Combine categorical feature and numerical feature to generate a new one. - Categorical Operations: Combine two categorical features to generate a new feature. For example, you 885 can infer a condition to make a binary feature, indicating whether it follows the condition. 886 887 You also have some new example features generated with these modules. 888 889 **Example Features:** Index | Feature_name | Feature_desc 890 1 | workclass_gender | Combine workclass and gender to see if certain workclasses have gender-specific 891 income differences 892 2 | education_marital | Combine education and marital-status to see if higher education levels are 893 correlated with higher income levels 894 3 | hours_per_week_squared | Square the hours-per-week column to see if there is a non-linear relationship with income 895 4 | age_log | Take the logarithm of age to normalize the distribution 896 5 | capital_net | Subtract capital-loss from capital-gain 897 You must write new feature that is different from all above examples features with respect to 899 both names and descriptions. 900 Format of response for 5 new columns: 901 902 Thought 1: [Any reasons based on examples above why the following new feature would be helpful for 903 the task] 904 New feature 1: [Type of operation] | New_column_name | One line detailed pseudo code for generating columns 905 906 Thought 5: ... 907 New feature 5: ... 908 909 Answer: 910 911 Thought 1: 912 913 914

Figure 7: Full prompt example for feature discovery in the Adult dataset (subsequent query with diversity enforcement). Text colored blue represents the diversity enforcement component.

A.3 EXAMPLE GENERATION RESULT FROM FEATURE DISCOVERY PROMPT

Figure 8 is an example of the outcome when a prompt for feature discovery is input into the LLM. A rationale and a simple description for the implementation of each new feature are generated.

Thought 1: Based on the examples, we can see that capital-gain and capital-loss could be important factors in determining if a person earns more than 50000 dollars per year. People with higher capital gains are more likely to earn more. New feature 1: Categorical Operations | HighEducation | If education is higher than HS-grad, set HighEducation to 1, else 0. Thought 2: The occupation of an individual can also be indicative of their income level. New feature 2: Categorical Operations | HighIncomeOccupation | If the occupation is Exec-managerial or Prof-specialty, set HighIncomeOccupation to 1, else 0. Thought 3: The capital gain and capital loss could be combined to create a net cap-ital feature. New feature 3: Numerical Operations | NetCapital | NetCapital = capital-gain - capital-loss Thought 4: The combination of age and hours-per-week could provide insights into the individual's work ethic and experience. New feature 4: Numerical Operations | WorkEthic | WorkEthic = age * hours-per-week Thought 5: The combination of marital status and relationship could indicate the in-dividual's family situation, which might affect their income. New feature 5: Mixed-type Operations | FamilyStatus | If marital-status is Married-civ-spouse and relationship is Husband or Wife, set FamilyStatus to 1, else 0. Figure 8: Example generation result by LLM from feature discovery prompt in the Adult dataset.

972 A.4 FULL PROMPT EXAMPLE FOR PARSING FEATURES

Figure 9 shows the prompt used to create converter code for generating actual features based on the
new feature descriptions provided by the LLM during the feature discovery process. This prompt is
input into the LLM to extract Python-based conversion code for use.

977	
978	Provide me a python code for function given description below
979	r tovide me a python code for function, given description below.
980	Function name: column appender
981	Input: Dataframe df_input
982	Input Features:
983	- age: the age of an individual (numerical variable within range [17, 90])
984	
985	
986	Output: Dataframe df_output.
987	Function descriptions. Create a new deteformer df entert. Fack column in df entert
988	Function description: Create a new dataframe di_output. Each column in di_output
989	the function code well matches with its feature type (i.e. numerical categorical) and follows
990	the below guidelines
991	the below guidelines.
992	Code guidelines (Format - Type New column name One line pseudo code):
993	- Generate each column with one line code.
994	- If there is no sufficient information to generate code, ignore the corresponding column.
995	- Make sure that the resulting type column is either string (for categorical) or numeric (for
996	continuous).
997	
998	Column description:
999	New feature 1: Categorical Operations HighEducation If education is higher than HS-grad,
1000	set HighEducation to 1, else 0.
1001	New feature 2: Categorical Operations HighIncomeOccupation If the occupation
1002	is Exec-managerial or Prof-specialty set HighIncomeOccupation to 1 else 0
1003	is Exec manuferral of Fior specially, set fingline one occupation to 1, else o.
1004 1005	New feature 3: Numerical Operations NetCapital NetCapital = capital-gain - capital-loss
1006 1007	New feature 4: Numerical Operations WorkEthic WorkEthic = age * hours-per- week
1008 1009 1010	New feature 5: Mixed-type Operations FamilyStatus If marital-status is Married- civ-spouse and relationship is Husband or Wife, set FamilyStatus to 1, else 0.
1011 1012	Figure 9: Full prompt example for parsing features in the Adult datsaet.
1013	
1014	
1015	
1016	
1017	
1018	
1019	
1020	
1021	
1022	
1023	
1024	

1026 A.5 EXAMPLE GENERATION RESULT FOR PARSING FEATURES PROMPT

Figure 10 is an example of the LLM's output for the parsing features prompt. It provides a Pythoncode function for creating the discovered feature.

1031	<stort></stort>
1032	def column appender(df input):
1033	der comm_appender(dr_mput).
1034	df output = df input.copy()
1035	df_output['HighEducation'] = df_output['education'].apply(lambda x: 1 if x !=
1036	'HS-grad' else 0)
1037	df_output['HighIncomeOccupation'] = df_output['occupation'].apply(lambda x: 1 if x
1038	in ['Exec-managerial', 'Prof-specialty'] else 0)
1039	df_output['NetCapital'] = df_output['capital-gain'] - df_output['capital-loss']
1040	df_output['WorkEthic'] = df_output['age'] * df_output['hours-per-week']
1041	df_output['FamilyStatus'] = df_output.apply(lambda x: 1 if (x['marital-status'] ==
1042	'Married-civ-spouse' and x['relationship'] in ['Husband', 'Wife']) else 0, axis=1)
1043	return df_output
1044	< and >
1045	

Figure 10: Example generation results for parsing features prompt in the Adult dataset.

1080 B DATASET DETAILS

In our work, we use a total of 22 datasets. Each dataset is assigned to either a binary, multi-class classification, or regression task. The datasets were selected considering the size, variety, and types of features. Basic information of each dataset are shown in Table 3 below. Task objectives of entire datasets are listed in Appendix A.

Table 3	3:	Basic	inf	ormation	of	datasets	used	for	evaluation.
	•••			011110001011	~ -				

Data	# of samples	# of features (Categorical/Numerical)	Task
Adult	48842	14 (7/7)	Binary classification
Balance-scale	625	4 (0/4)	Multi-class classification
Bank	45211	16 (8/8)	Binary classification
Blood	748	4 (0/4)	Binary classification
Car	1728	6 (5/1)	Multi-class classification
Communities	1994	103 (1/102)	Multi-class classification
Credit-g	1000	20 (12/8)	Binary classification
Diabetes	768	8 (0/8)	Binary classification
Eucalyptus	736	19 (5/14)	Multi-class classification
Forest-fires	517	12 (2/10)	Regression
Heart	918	11 (4/7)	Binary classification
Junglechess	44819	6 (0/6)	Multi-class classification
Myocardial	1700	111 (94/17)	Binary classification
Tic-tac-toe	958	9 (9/0)	Binary-classification
Vehicle	846	18 (0/18)	Multi-class classification
Bike	17379	12 (3/9)	Regression
Crab	3893	8 (1/7)	Regression
Housing	20640	9 (1/8)	Regression
Insurance	1338	6 (3/3)	Regression
Wine	6497	12 (1/11)	Regression
Sequence-type	250	5 (0/5)	Multi-class classification
Solution-mix	300	8 (0/8)	Regression

1134 С **IMPLEMENTATION DETAILS** 1135

1136 C.1 TST-LLM DETAILS 1137

1138 This section provides additional implementation details of our model. In the feature discovery process of TST-LLM, we use the GPT-3.5 model as the LLM backbone. Meta-information such as feature 1139 names and descriptions were included in the prompt. For categorical features, a list of categories 1140 for each feature was added, and for numerical features, the min-max value statistics were included. 1141 During LLM generation, the temperature was set to 0.5 and the top-p value was set to the API's default 1142 of 1. The discovery process generated five features per trial, with the number of trials set at 40. The 1143 number of serialized samples included in the prompt was set to a maximum of 20, as allowed by the 1144 prompt limit. When the number of features in a dataset exceeded 100 (e.g., communities, myocardial), 1145 and the prompt limit was reached, we resolved this by selecting a random 10 columns per query. 1146 Over 40 trials, we ensured that all features were used at least once in the feature discovery process. 1147

After the LLM completed feature discovery, a feature set satisfying the minimum redundancy 1148 between the original data was selected for representation learning. The number of selected features 1149 M was set to 20. Refer to Algorithm 1 below for the feature selection algorithm. The encoder 1150 structure for representation learning was consistent with the baselines, configured as a 2-layer MLP 1151 with 1024 dimensions. The projection head consisted of a single linear layer, projecting 1024 to 1152 128 dimensions. Training utilized the Adam optimizer with a learning rate of 1e-4, a batch size of 1153 128, and 1000 training iterations.

1154 1155

Algorithm 1: Algorithm for feature selection with minimum redundancy.

1156		Argorithm 1. Argorithm for readile selection with minimum redundancy.
1157		Input :Initial feature set $\hat{\mathcal{Y}}_{init}$, number of features to select M , original dataset \mathcal{D} .
1158	1	$\hat{\mathcal{V}} \leftarrow \emptyset$
1159	2	$\hat{\mathcal{Y}}_{\text{end}} \leftarrow \{\hat{u} \mid \text{Entrony}(\hat{u}) > t_{\text{end}} \mid \hat{u} \in \hat{\mathcal{Y}}_{\text{end}}\}$
1160	2	while $ \hat{\mathcal{Y}} < M$ do
1161	3 4	$ \Phi \leftarrow \emptyset$
1162	5	for $\hat{y} \in \hat{y}_{charact}$ do
1163	6	$\phi_{u} \leftarrow \max(\operatorname{CramersV}(\mathcal{D}, \hat{y}));$ // Compute redundancy of the feature
1164	7	$ \left \begin{array}{c} \Phi \leftarrow \Phi \cup \{(\phi_y, \hat{y})\} \end{array} \right $
1165	8	end
1166		/* Select features with minimum redundancy */
1167	9	$\hat{\mathcal{Y}}_{\text{selected}} \leftarrow \{ \hat{y} \mid \phi_y = \min_{\phi_y}(\Phi), (\phi_y, \hat{y}) \in \Phi \}$
1168	10	$\hat{\mathcal{Y}}_{ ext{filtered}} \leftarrow \hat{\mathcal{Y}}_{ ext{filtered}} - \hat{\mathcal{Y}}_{ ext{selected}}$
1169	11	$\mathcal{D} \leftarrow \mathcal{D} \cup \hat{\mathcal{Y}}_{ ext{selected}}$
1170	12	$\hat{\mathcal{Y}} \leftarrow \hat{\mathcal{Y}} \cup \hat{\mathcal{Y}}_{ ext{selected}}$
1171	13	end
1172		
1173		
1174		
1175		
1176		
1177		
1178		
1179		
1180		
1181		
1182		
1183		
1184		
1185		
1186		
1187		

1188 C.2 BASELINE DETAILS

This section describes the implementation details of the baselines. While the implementation of the baselines followed the original works of the respective papers, the encoder used to extract representations was configured uniformly for a fair comparison (i.e., a 2-layer MLP with 1024 hidden dimensions). Different decoder and projector networks were used according to each methodology.

For Autoencoder baseline, the decoder was the same 2-layer MLP with 1024 hidden dimensions as the encoder. For Siamese network-based methodologies (e.g., SimSiam, SCARF, STAB), a 2-layer MLP with 256 hidden dimensions was used as the projector, and for SimSiam, the predictor consisted of a single linear layer. STUNT, which uses prototype-based learning, does not have a separate decoder. For LFR, a single linear layer predictor and a 2-layer ReLU network with 256 hidden dimensions were used as the random data projector.

For all baselines, we referred to the following links for the implementation²³.

²https://github.com/layer6ai-labs/lfr ³https://github.com/jaehyun513/STUNT

¹²⁴² D COMPUTATIONAL COMPLEXITY

1243 1244 COMPUTATIONAL COMPLEXITY

In this section, we compare the computational time required for model training. The comparison was conducted on the Adult dataset using a single A100 GPU. For our model, the computation time includes the entire process of feature discovery and selection from the LLM, as well as training.
Table 4 reports the total time spent for each method. We found that our model has a computational time complexity comparable to other baselines.

Table 4: Computational time complexity analysis of self-supervised representation learning methods.
The total time spent (in seconds) and the ratio compared to our model are reported for each method.

1252			
1253	Model	Time spent (second)	Time spent (ratio)
1254	Autoencoder	520.4	1.08
1255	SimSiam	350.7	0.73
1256	SCARF	479.6	1.00
1257	STAB	208.7	0.43
1258	STUNT	608.8	1.27
1259	LFR	470.3	0.98
1260	TST-LLM	481.2	1.00
1261			
1262			
1263			
1264			
1265			
1266			
1267			
1268			
1269			
1270			
1271			
1272			
1273			
1274			
1275			
1276			
1277			
1278			
1279			
1280			
1281			
1282			
1283			
1284			
1285			
1286			
1287			
1288			
1289			
1290			
1291			
1292			
1293			
1294			
1295			

1296 E LEARNING WITH OTHER OBJECTIVES

1298 In our framework, we utilize supervised contrastive learning to integrate information from LLM-1299 discovered features into embeddings, although it is not the only available approach. Therefore, in 1300 this section, we compare the performance of our framework using different loss objectives with a 1301 linear model. Figure 11 compares the performance of self-supervised baselines and TST-LLM against 1302 each other using win matrices, while our framework uses different training objectives including supervised contrastive learning (Figure 11a), CLIP (Radford et al., 2021) (Figure 11b), reconstruction 1303 (Figure 11c), and cross-entropy (Figure 11d). Our framework consistently outperforms other self-1304 supervised baselines, irrespective of the training objectives used. 1305



Figure 11: Win matrices comparing self-supervised tabular learning methods against each other, while our framework uses different training objectives including (a) Supervised Contrastive Learning, (b) CLIP, (c) Reconstruction, and (d) Cross-Entropy. Self-supervised tabular learning methods are aligned on the x-axis and the y-axis while the numbers represent the winning ratio of the x-axis model against the y-axis model. Full results are in the Appendix H.7.

- 1339 1340
- 1341
- 1343
- 1343
- 1345
- 1346
- 1347
- 1348
- 1349

1350 F IMPACT OF THE NUMBER OF TRIALS IN FEATURE DISCOVERY

In this section, we analyze the impact of the number of trials on downstream task performance when performing feature discovery through an LLM. In our current model setting, five new features are discovered per trial, and a total of 40 trials are made to obtain the feature set. Figure 12 below measures the performance change ratio compared to the current model as the number of trials is varied to 5, 10, 20, and 30. The results indicate that with 10 or more trials, stable performance is achieved across multiple tasks.



Figure 12: Effect of the number of trials in feature discovery on the performance of downstream tasks.



1402 1403

1358 1359

1360

1361 1362

1404 G QUALITATIVE ANALYSIS

To verify whether the features discovered by the LLM align with the task definition, we selected and examined the top three discovered features for each dataset using our selection strategy (see Table 5).
We observed that the discovered features somewhat intuitively align with the downstream task.

- 1409
- 1410

Table 5: Top-3 discovered features from our selection strategy for each dataset.

Data Adult Balance-scale Bank Blood Car	Top-3 discovered features age * hours-per-week, educational-num / age, educational-num * age abs(left-weight - right-weight), abs(left-weight + left-distance - right-weight - right distance), (left-weight - left-distance)**2 - (right-weight - right-distance)**2 duration * campaign, balance * duration, duration / day (Recency ** 0.5) * Frequency / Time, 1 / (1 + np.exp(Recency - Time)), (Time - Recency) / Frequency
Adult Balance-scale Bank Blood Car	age * hours-per-week, educational-num / age, educational-num * age abs(left-weight - right-weight), abs(left-weight + left-distance - right-weight - right distance), (left-weight - left-distance)**2 - (right-weight - right-distance)**2 duration * campaign, balance * duration, duration / day (Recency ** 0.5) * Frequency / Time, 1 / (1 + np.exp(Recency - Time)), (Time - Recency) / Frequency
Balance-scale Bank Blood Car	abs(left-weight - right-weight), abs(left-weight + left-distance), (left-weight - left-distance)**2 - (right-weight - right-distance)**2 duration * campaign, balance * duration, duration / day (Recency ** 0.5) * Frequency / Time, 1 / (1 + np.exp(Recency - Time)), (Time - Recency) / Frequency
Bank Blood Car	duration * campaign, balance * duration, duration / day (Recency ** 0.5) * Frequency / Time. 1 / (1 + np.exp(Recency - Time)), (Time - Recency) / Frequency
Blood	(Recency ** 0.5) * Frequency / Time. 1 / (1 + np.exp(Recency - Time)), (Time - Recency) / Frequency
Car	
	<pre>maint.map({'low': 1, 'medium': 2, 'high': 3, 'very high': 4}) + doors.map({'5more': 5, '4': 4, '3': 3, '2': 2}), buying.map({'high':3, 'low':1, 'medium':2, 'very high':4}) + maint.map({'high':3, 'low':1, 'medium':2, 'very high':4}) + doors.map({'5more':4, '4':3, '2':1, '3':2}) + persons.map({'more':4, '2':1, '4':3}) + lug_boot.map({'med':2, 'big':3, 'small':1}) + safety.map({'med':2, 'low':1, 'high':3}), (maint + safety)/2</pre>
Communities	PctEmplManu * HousVacant, MedRentPctHousInc * pctWWage, agePct12t21 * NumInShelters
Credit-g	duration / age, age * duration, age / duration
Diabetes	Glucose / Age, Pregnancies * DiabetesPedigreeFunction, DiabetesPedigreeFunction.map(DiabetesPedigreeFunction.value_counts())
Eucalyptus	(Surv + Vig) * Ht, Stem_Fm - Brnch_Fm, Crown_Fm - Brnch_Fm
Forest-fires	temp * RH, wind * temp, FFMC + DMC + DC + ISI
Heart	Age.corr(MaxHR), RestingBP * MaxHR, abs(RestingBP - MaxHR)
Junglechess	(white_piece0_file * white_piece0_rank) / (black_piece0_file * black_piece0_rank), white_piece0_file * white_piece0_rank, groupby(('white_piece0_file', 'white_piece0_rank', 'black_piece0_file', 'black_piece0_rank)['white_piece0_file'].transform('count'
Myocardial	log(AST_BLOOD), L_BLOOD * ROE, L_BLOOD.value_counts()[L_BLOOD].values
Tic-tac-toe	apply(lambda x: (x['top-left-square'] == 'x') + (x['middle-middle-square'] == 'x') + (x['bottom-right-square'] == 'x')), apply(lambda x: (x['bottom-left-square'] == 'o') + (x['bottom-middle-square'] == 'o') + (x['bottom-right-square'] == 'o')), apply(lambda x: [x['top-left-square'], x['top-right-square'], x['bottom-left-square'], x['bottom-right-square'].count('o')),
Vehicle	COMPACTNESS / CIRCULARITY, SCALED_RADIUS_OF_GYRATION / RADIUS_RATIO, PR.AXIS_RECTANGULARITY / CIRCULARITY
Bike Crab	abs(temp - hum), abs(temp - atemp), hr * mnth Shell Weight / (Weight - Shucked Weight - Viscera Weight), Shucked Weight / Viscera Weight, Weight.value_counts()
Housing Insurance	population / households, total_bedrooms / households, median_income / population age * bmi, abs(age - bmi), age / (children + 1)
Wine	sulphates - volatile acidity, citric acid / residual sugar, fixed acidity + alcohol
Sequence-type	Number2 / Number1 - Number3 / Number2, ['Number1', 'Number2', 'Number3', 'Number4', 'Number5'].sum(axis=1) % 2, (Number2 / Number1 + Number3 / Number2 + Number4 / Number3 + Number5 / Number4).cumsum()
Solution-mix	Solution_1_volume * Solution_1_concentration + Solution_2_volume * Solution_2_concentration + Solution_3_volume * Solution_3_concentration + Solution_4_volume * Solution_4_concentration, abs(Solution_1_concentration - Solution_2_concentration) + abs(Solution_2_concentration - Solution_3_concentration) + abs(Solution_3_concentration - Solution_4_concentration), np.log(Solution_1_concentration + Solution_2_concentration + Solution_3_concentration) / (Solution_1_volume + Solution_2_volume + Solution_3_volume + Solution_4_volume))

FULL RESULTS Η

In this section, we present full results of our experiments.

H.1 EVALUATION WITH LINEAR MODEL

Table 6: Evaluation results of self-supervised models on linear model, showing (a) AUC across 15 datasets for classification and (b) RMSE across 7 datasets for regression. Best performances are bolded, and our framework's performances, when second-best, are underlined.

(a) Classification (AUC)

Dataset	Raw Data	AutoEncoder	SimSiam	SCARF	STAB	STUNT	LFR	Ou
Adult	90.75±0.17	91.07±0.20	89.01±0.24	4 90.90±0.17	90.55±0.24	91.06±0.20	91.29±0.19	91.32
Balance-scale	97.24±1.11	99.58±0.37	99.37±0.4	5 99.44±0.39	97.66±1.46	93.10±2.50	99.28±0.39	99.51:
Bank	$90.48 {\pm} 0.18$	$91.14 {\pm} 0.07$	87.32±0.1	5 91.73±0.04	$90.06 {\pm} 0.24$	91.10±0.38	$91.65 {\pm} 0.17$	92.08
Blood	75.15 ± 3.21	$74.98 {\pm} 3.52$	75.18±4.4	2 73.92±4.04	74.75 ± 3.24	$74.39 {\pm} 4.83$	$73.88 {\pm} 3.11$	74.85
Car	$98.95 {\pm} 0.30$	99.60 ± 0.23	97.95±0.42	2 99.50±0.31	99.25±0.43	97.96 ± 0.28	99.91±0.04	<u>99.73</u>
Communities	84.31±1.23	83.01 ± 0.74	83.56±0.8	5 83.86±1.51	84.39 ± 0.70	85.09±1.15	81.45 ± 1.12	85.25
Credit-g	77.89 ± 6.44	77.60 ± 5.26	77.69 ± 6.2	7 77.12 \pm 5.46	77.26 ± 3.18	75.94 ± 4.51	75.04 ± 6.32	78.38
Diabetes	83.07±4.74	81.64 ± 6.16	82.73±5.5	8 81.96±5.97	80.38 ± 5.22	82.21±3.04	81.43 ± 7.38	82.56
Eucalyptus	91.64±1.10	90.85 ± 1.33	90.50 ± 1.8	90.44 ± 1.23	89.66 ± 1.97	85.61 ± 1.51	89.85 ± 0.64	91.34
Heart	93.10 ± 2.12	92.79 ± 1.60	93.07±2.3	3 93.15±1.58	93.15 ± 2.52	92.38 ± 2.70	92.60 ± 2.16	93.45
Junglechess	80.61 ± 0.33	89.89 ± 0.49	86.92 ± 0.7	88.45 ± 0.70	92.10 ± 0.47	91.62 ± 0.44	92.93 ± 0.42	93.43
Myocardial	61.20 ± 5.13	60.90 ± 4.97	66.11 ± 4.0	$5 60.43 \pm 3.35$	59.29 ± 3.86	63.27±4.35	62.06 ± 3.38	63.64
Sequence-type	92.11 ± 2.03	96.37±0.75	96.34±1.1	7 97.36±0.91	97.16±1.48	92.40±1.15	97.41±0.63	96.44
Tic-tac-toe	99.31±0.60	99.84±0.08	98.28±1.3	5 99.00±0.67	95.93±1.87	94.07±3.14	99.80±0.15	99.52
Vehicle	94.82 ± 0.50	96.16±0.83	92.37 ± 1.3	996.02 ± 0.52	95.32±0.49	93.55±1.07	96.32±0.38	<u>96.22</u>
			(0) K	egression (KN	13E)			
	Dataset	Raw Da	ta A	utoEncoder	SimSian	n S	CARF	
	Bike	142.36±1	.58 1	26.90±1.02	121.59±1.	59 111	.67±2.08	
	Crab	2.21 ± 0.0	05	$2.12{\pm}0.03$	$2.12{\pm}0.0$	4 2.1	12±0.03	
	Forest-fires	75.07±35	.28 8	$1.21{\pm}29.08$	82.01±27.	45 82.8	37 ± 28.15	
	Housing	69132.79±4	89.67 581	55.43±619.72	59159.48±6	2.79 56941	.63±519.79	
	Insurance	5930.14±27	73.29 464	41.78±220.29	4666.29±25	2.95 4657.	87 ± 174.01	
	Solution-mix	0.07 ± 0.0	00	$0.03 {\pm} 0.00$	$0.03 {\pm} 0.0$	0.0	03 ± 0.00	
	Wine	0.73±0.0	01	$0.69 {\pm} 0.00$	$0.69{\pm}0.0$	0 0.6	69±0.01	
	Dataset	STAB		STUNT	LFR		Ours	
	Bike	126.42+1	91 1	15 98+2 18	121 02+2	43 111	46+1 72	
	Crab	2.15 ± 0.0)2	2.13 ± 0.04	2.16±0.02	2 2.1	13±0.03	
	Forest-fires	80.13±30	.24 7	7.57±32.56	83.84±26.	81 83.1	19±22.47	
	Housing	60071.46 ± 29 4787 10+17	97.95 561: 0.53 500	51.34 ± 352.44	58064.28 ± 31 4833.99 ± 29	.2.73 56069 3.63 4578	.83±406.99 14+149 83	
	Solution-mix	0.03±0.0	0.55 509	0.07 ± 0.00	0.02±0.0	0 0.0	02 ± 0.01	
	Wine	0.71±0.0	00	0.67±0.00	$0.69{\pm}0.0$	1 0.	57±0.00	

1512 H.2 EVALUATION WITH NON-PARAMETRIC CLASSIFIER1513

Table 7: Evaluation results of self-supervised models on Non-parametric classifier with (a) 3 and (b)
 5 clusters, showing AUC across 16 datasets for classification. Best performances are bolded, and our
 framework's performances, when second-best, are underlined.

1	5	1	7
1	5	1	8

1563 1564 1565 (a) 3 Clusters

	Raw Data	AutoEncoder	SimSiam	SCARF	STAB	STUNT	LFR	Ours
Adult	82.31±0.18	81.53±0.16	80.51±0.41	81.90±0.09	82.21±0.24	82.20±0.30	82.28±0.29	81.90±0.46
Balance-scale	79.47±2.44	78.40±1.60	82.67±1.85	78.67±0.46	79.73±3.23	79.20±1.60	78.13±2.44	79.73±0.46
Bank Blood	89.09 ± 0.14 72 44 + 4 44	89.19 ± 0.03 71.33+3.46	88.41 ± 0.25 71 33+4 67	88.85 ± 0.13 69.78 + 5.00	89.11 ± 0.04 71.33 \pm 4.67	88.96 ± 0.11 74 67 + 3 71	89.16 ± 0.11 72 00 ± 4.16	89.36±0.26 74.33+3.06
Car	87.09±2.09	86.42 ± 2.37	77.17±1.16	79.00 ± 0.73	82.01±1.59	82.85 ± 1.97	72.00 ± 4.10 79.77 ± 2.08	$\frac{74.55\pm5.00}{89.40\pm2.73}$
Communities	$63.07{\pm}2.25$	$62.32{\pm}2.43$	$58.23 {\pm} 1.61$	$61.99{\pm}1.63$	$61.57{\pm}0.95$	$62.24{\pm}1.04$	$62.91{\pm}0.90$	$62.16{\pm}2.39$
Credit-g	73.00 ± 1.50	73.33±0.76	69.00 ± 4.44	71.83 ± 3.33	72.00 ± 3.04	71.67 ± 4.65	71.50 ± 1.00	71.83 ± 1.76
Eucalyptus	73.10 ± 4.90 59 23+3 19	72.94 ± 0.37 53 60+3 96	57.88 ± 2.73	52.03 ± 4.05	74.24 ± 1.35 56 08+4 22	74.03 ± 4.00 52.25+2.06	72.31 ± 4.12 58 23+2 06	60.59 ± 5.12
Heart	84.60±1.66	84.96±1.91	84.42±1.13	85.33±0.54	85.51±1.75	85.05±2.57	84.60±2.45	84.70±1.09
Junglechess	75.08±0.54	74.35±0.27	77.40±0.14	72.34±0.22	74.84±0.64	73.65±0.47	73.87±0.38	$\frac{75.35\pm0.52}{72.22\pm0.74}$
Myocardial	74.40 ± 5.63 90.00+2.00	73.91 ± 2.51 91.33 ±1.15	73.43 ± 2.93 86.00 ±3.46	75.12 ± 1.11 93 33+1 15	73.43 ± 1.11 90.67 ± 2.31	71.01 ± 3.32 92.00+2.00	74.64±4.35	72.22 ± 2.74 91.33+2.31
Tic-tac-toe	91.15 ± 0.52	82.29±0.90	85.07±2.41	72.40 ± 2.71	84.90±1.56	96.70±1.59	77.95±2.46	93.92±1.59
Vehicle	$69.80{\pm}2.96$	75.10±3.55	$59.80{\pm}2.96$	$68.82{\pm}5.23$	$68.80{\pm}1.80$	$67.59{\pm}3.67$	74.31±4.34	74.51 ± 2.65
			(b)) 5 Clusters				
Dataset	Raw Data	AutoEncoder	SimSiam	SCARF	STAB	STUNT	LFR	Ours
Adult	83.17±0.19	82.48±0.12	81.47±0.31	82.60 ± 0.11	83.17±0.11	82.93±0.35	83.15±0.37	83.22±0.32
ыalance-scale Bank	82.40 ± 2.88 89.40+0.32	82.40±2.12 89.48+0.14	88.85+0.22	δ5.20±2.40 89.13+0.16	81.00 ± 3.20 89.54 ± 0.04	85.07±2.81 89.37+0.26	81.33 ± 3.33 89.47 ± 0.27	81.33±0.92 89.25+0.06
Blood	74.67±4.16	74.44 ± 4.07	73.56 ± 4.73	74.44 ± 2.78	74.67±2.91	74.22 ± 4.34	72.89 ± 3.67	73.78±4.07
Car	89.69±0.83	84.90±1.01	78.13±1.77	85.07±1.64	88.54±2.53	88.54±2.67	84.78±1.92	93.77±2.29
Communities Credit-g	65.58±1.16 72 33+2 57	64.55±1.24 74 83+0 29	61.32 ± 0.14 71 17+4 54	63.16 ± 1.57 73 50 ± 3.61	63.41 ± 1.96 71.83+2.93	65.25 ± 1.67 71 33+1 15	64.91 ± 1.15 71.67 + 1.53	65.58 ±1.43 73.17+3.88
Diabetes	73.38 ± 3.25	72.51±4.32	73.38 ± 5.15	73.16±2.46	72.94 ± 3.33	74.03±3.62	73.38 ± 3.90	71.65 ± 5.52
Eucalyptus	59.46±4.87	54.95±3.96	58.33±2.56	53.60±3.47	54.28±5.46	52.70±3.10	58.33±3.96	63.16±3.10
Heart Junglechess	85.69 ± 1.91 75 20 ± 0.42	86.23 ± 0.31 75 28 ± 0.40	85.69±1.57 78 93±0 57	86.41 ± 1.44 74.04 ± 0.54	84.96 ± 0.63 76.09 \pm 0.61	85.69 ± 0.83 75.57 \pm 0.33	86.59±1.66	84.06 ± 1.37 75 80 ± 0.43
Myocardial	75.60 ± 2.74	75.85 ± 1.82	73.91±0.72	76.33 ± 1.11	76.57±1.82	74.40 ± 2.54	76.09 ± 2.90	76.12 ± 2.33
Sequence-type	$91.33 {\pm} 3.06$	$91.33{\pm}1.15$	$86.00{\pm}5.29$	93.33±2.31	90.00 ± 3.46	90.00 ± 4.00	93.33±2.31	$92.67 {\pm} 3.06$
Tic-tac-toe	94.10 ± 0.80 72 75+1 70	84.38 ± 1.80 75.49 ± 0.90	87.33 ± 1.97 60 39 ± 0.34	77.26 ± 2.46 70.39 ± 3.59	90.45 ± 2.87 72 75 ± 0.34	97.74±0.60 71.57±2.23	82.12 ± 1.08 74 71 + 3 11	$\frac{95.31\pm1.88}{76.47\pm3.53}$
vennene	12002100	7010)±0100	00107±0101	10109 20109	/2//0 ±010 1	/110/±2120	,	

H.3 ABLATION STUDY WITH LINEAR MODEL

Table 8: Evaluation results of ablation studies on linear model, showing (a) AUC across 15 datasets for classification and (b) RMSE across 7 datasets for regression. Best performances are bolded, and our framework's performances, when second-best, are underlined.

(a) Classification (AUC)

Dataset	Top-1 selection	Random-1 selection	Random feature discov	ery Without learning	Without feature selection	on Ours
Adult	91.27±0.11	91.27±0.13	91.38±0.11	91.65±0.54	91.36±0.13	91.32±0.12
Balance-scale	99.46±0.37	99.53±0.29	99.54±0.26	99.96±0.07	99.64±0.26	99.51±0.33
Bank	$92.10 {\pm} 0.12$	91.96 ± 0.20	91.99 ± 0.24	$89.83 {\pm} 0.65$	91.91±0.19	92.08 ± 0.18
Blood	74.72 ± 2.85	$74.89 {\pm} 2.83$	74.85 ± 2.85	73.12 ± 4.94	$75.26{\pm}2.67$	74.85 ± 2.89
Car	99.65 ± 0.22	$99.48 {\pm} 0.40$	99.68±0.16	98.05±1.75	99.81±0.08	99.73±0.18
Communities	85.37±1.28	84.74 ± 0.44	85.36 ± 0.73	84.19 ± 1.46	85.59±1.28	85.25 ± 0.67
Credit-g	78.23 ± 5.14	77.97 ± 5.89	77.53 ± 5.48	$74.48 {\pm} 4.54$	78.43±4.46	78.38 ± 4.85
Diabetes	82.20 ± 5.36	82.23 ± 5.19	82.60±4.81	84.33 ± 3.89	82.19 ± 5.20	82.56 ± 5.12
Eucalyptus	91.16±0.84	90.92 ± 1.26	91.15±0.89	$88.94{\pm}0.80$	91.41±1.13	91.34 ± 0.99
Heart	93.56±1.44	93.46 ± 1.28	93.25±1.47	92.50±1.63	93.47±1.48	93.45±1.60
Junglechess	93.37±0.28	92.07 ± 0.29	93.37±0.21	93.57±1.54	93.43±0.39	93.43 ± 0.31
Myocardial	62.10 ± 2.45	63.05 ± 1.51	62.46 ± 2.78	$63.98{\pm}2.70$	62.23 ± 3.70	63.64 ± 3.08
Sequence-type	96.58 ± 0.90	$96.80{\pm}1.02$	96.45 ± 1.01	83.33 ± 28.87	96.47 ± 1.10	96.44 ± 1.01
Tic-tac-toe	99.58±0.23	99.24 ± 0.42	98.95 ± 0.53	99.95±0.06	99.47±0.33	99.52 ± 0.51
Vehicle	$95.95 {\pm} 0.51$	$95.94{\pm}0.47$	$96.00 {\pm} 0.56$	$94.69 {\pm} 0.30$	$96.08 {\pm} 0.53$	96.22±0.28
			(b) Regression (R	MSE)		
Dataset	Top-1 selection	Random-1 selection	Random feature discovery	Without learning	Without feature selection	Ours
Bike	112.70±1.83	111.27±1.83	111.09±1.75	740.74±657.86	111.68±2.18	111.46±1.72
Crab	$2.12{\pm}0.03$	$2.12{\pm}0.01$	2.13 ± 0.02	$6.94{\pm}6.20$	2.13 ± 0.02	2.13±0.03
Forest-fires	81.61±22.83	81.41±23.53	83.55±23.47	87.11±32.62	86.03±24.27	83.19±22.47
Housing	56162.70±247.14	56231.24±334.36	55984.82 ± 323.18	65521.90±9492.64	55864.26±90.89	56069.83±406.99
Insurance	4630.87±139.19	4567.01±116.92	4587.91 ± 118.33	5880.78 ± 1411.86	4582.00±179.75	4578.14 ± 149.83
Solution-mix Wine	0.02 ± 0.01 0.68 ± 0.00	0.02 ± 0.00 0.68 ± 0.00	0.02 ± 0.00 0.68 ± 0.00	0.01±0.00 0.90±0.26	0.02 ± 0.00 0.68 ± 0.01	$\frac{0.02\pm0.01}{0.67\pm0.00}$

1620 H.4 ABLATION STUDY WITH NON-PARAMETRIC CLASSIFIER

1622Table 9: Evaluation results of ablation studies on non-parametric classifier with (a) 3 and (b) 51623clusters across 16 datasets for classification. Best performances are bolded, and our framework's1624performances, when second-best, are underlined.

1	6	2	5
1	6	2	6

(a) 3 Clusters

Dataset	Top-1 selection	Random-1 selection	Random feature discovery	Without learning	Without feature selection	Ours
Adult	$81.87{\pm}0.41$	$82.11 {\pm} 0.08$	82.64±0.36	83.07±0.23	83.04±0.21	81.90±0.46
Balance-scale	73.07±5.45	79.20±1.39	80.00±1.60	86.13±0.46	80.27±0.92	79.73±0.46
Bank	89.08±0.05	89.04±0.10	88.91±0.36	88.76±0.29	88.71±0.19	89.36±0.26
lood	$\frac{1.33\pm1.0}{82.56\pm0.62}$	72.44 ± 5.05	/3.33±3.06 82.66±6.05	76.00±2.00 81 79±5 60	/2.6/±3.53 83.82±3.33	$\frac{14.33\pm3.06}{89.40\pm2.72}$
aı ommunities	63.32 ± 2.27	69.02 ± 2.57 64.24 ± 2.04	62.00 ± 0.93 62.74 ± 0.29	63.24 ± 1.09	63.62±3.55 64.75+1.16	62.16 ± 2.73
redit-g	70.00 ± 2.60	71.33 ± 2.57	72.33±1.76	72.00 ± 5.63	71.00 ± 2.65	71.83±1.76
abetes	$72.94{\pm}5.52$	72.29 ± 4.32	$71.94{\pm}4.92$	73.38±5.15	$72.94{\pm}6.14$	$72.73 {\pm} 4.90$
ıcalyptus	56.53 ± 3.96	$59.68 {\pm} 2.73$	58.11±1.17	$56.53 {\pm} 2.81$	$62.39{\pm}4.50$	60.59 ± 5.12
art	85.69±1.75	84.24 ± 2.72	84.05±1.57	86.78±0.83	84.06±2.57	84.70±1.09
glechess	/5.46±1.21	75.45 ± 0.75	74.91±0.63	75.76±2.05	74.15±0.52	75.35 ± 0.52
uence-type	75.12 ± 4.25 90.00 ± 3.46	75.19 ± 5.52 80 33+2 31	73.07 ± 0.84 91.33+2.31	72.40 ± 4.55 73.33 ±20.14	74.13 ± 2.55 92 00+3 46	72.22 ± 2.74 01 33+2 31
c-tac-toe	90.80 ± 1.20	91.49 ± 1.50	81.25±5.29	91.15 ± 12.18	95.49±1.31	$\frac{91.95\pm2.91}{93.92\pm1.59}$
ehicle	71.76±5.79	70.59 ± 5.79	72.75±4.42	74.12±0.59	70.78 ± 1.48	74.51±2.65
			(b) 5 Clusters			
taset	Top-1 selection	Random-1 selection	Random feature discovery	Without learning	Without feature selection	Ours
ult	82.64±0.15	82.90±0.13	83.43±0.35	84.03±0.11	83.83±0.22	83.22±0.32
lance-scale	$81.33 {\pm} 3.23$	$82.93 {\pm} 2.44$	$82.93 {\pm} 1.22$	86.13±2.44	$81.07 {\pm} 2.01$	$81.33{\pm}0.92$
nk	89.51±0.16	89.37±0.13	89.21 ± 0.20	89.41 ± 0.18	89.25 ± 0.08	89.25±0.06
lood	73.78±2.14	75.33 ± 3.06	73.56±3.79	76.00±3.06	72.67±2.91	73.78±4.07
ſ mmunitias	88.82±6.79	91.04 ± 2.00 65.33 ±1.42	88.54±3.47	85.62 ± 5.43	89.69±0.60	95.77±2.29
redit_g	7317 ± 1.13	733 + 736	04.41 ± 1.23 73.00+2.60	04.24 ± 1.13 71 33+5 53	04.05 ± 2.05 71.67 ±3.62	03.30±1.43
abetes	73.59 ± 5.67	73.59±3.33	74.24±4.92	72.29 ± 4.61	73.16 ± 4.70	$\frac{75.17\pm5.08}{71.65\pm5.52}$
ucalyptus	58.56±7.29	59.46 ± 5.41	56.76±4.11	58.11±3.76	62.39±6.39	63.16±3.10
art	$85.69 {\pm} 1.91$	$85.14{\pm}1.66$	$87.14{\pm}1.13$	$85.87 {\pm} 0.54$	$84.24{\pm}1.44$	$84.06 {\pm} 1.37$
nglechess	76.51±1.35	75.30±0.31	75.73 ± 0.60	76.94±1.88	75.05±0.41	75.80 ± 0.43
yocardial	75.85 ± 1.67	73.91 ± 0.72	75.36 ± 1.92	75.36±2.90	75.60 ± 1.11	76.12±2.33
quence-type	91.33±4.16 94.62±0.60	91.33 ± 4.10 92.53 ± 2.10	90.07±4.62 86.98±5.02	13.33 ± 29.48 91 15+11 30	92.00 ± 4.00 93.92+1.50	92.0/±5.06 95 31+1 89
hicle	72.75 ± 2.78	73.14 ± 3.02	74.51+1.36	70.59 ± 4.08	73.33+1.80	76.47+3.53
		.5.1.125.02	,	,010,911,00	75.55±1.00	

1674 H.5 INFORMATIVENESS OF DISCOVERED FEATURES

1680

Table 10: Analysis of the informativeness of features discovered via LLM. The average mutual information (MI) between features and the downstream task's labels is reported for each dataset. The increase ratio in MI when using discovered features compared to original features is also reported, along with standard deviations.

Data	Average MI in original features	Average MI in discovered features	Increase ratio (%)
Tic-tac-toe	0.010	0.076	646.6±1535.1
Solution-mix	0.064	0.386	507.3±1348.2
Balance-scale	0.082	0.178	117.3 ± 299.0
Wine	0.055	0.100	81.2±112.3
Bank	0.013	0.022	65.8±172.0
Blood	0.031	0.045	44.6 ± 95.2
Sequence-type	0.345	0.459	32.9 ± 86.9
Forest-fires	0.019	0.025	32.2±131.3
Bike	0.103	0.132	27.9 ± 155.0
Car	0.036	0.041	14.2 ± 157.5
Credit-g	0.009	0.009	8.8 ± 136.5
Insurance	0.364	0.395	$8.4{\pm}127.7$
Communities	0.089	0.095	7.1 ± 85.2
Vehicle	0.212	0.226	6.8 ± 56.5
Adult	0.031	0.032	5.3 ± 123.1
Junglechess	0.049	0.051	4.8 ± 73.8
Myocardial	0.007	0.007	3.1 ± 95.5
Diabetes	0.043	0.045	3.0 ± 71.9
Heart	0.067	0.063	-5.1 ± 81.5
Eucalyptus	0.1//	0.158	-10.9 ± 86.9
Uning	0.550	0.234	-27.5 ± 43.7

1728 H.6 HYPERPARAMETER ANALYSIS ON M

1730Table 11: Evaluation results with various hyperparameter M on linear model, showing (a) AUC1731across 15 datasets for classification and (b) RMSE across 7 datasets for regression. Best performances1732are bolded.

1733								
1734				(a) Cl	assification (AU	C)		
1735		Dataset	M =	10	M = 20	M = 30	M	= All
1736		Adult	01 22	0.16	01 22+0 12	01 21 + 0 14	01.2	7+0.16
1/3/		Auun Palanaa saala	$91.33 \pm 00.50 \pm$	0.10	91.32 ± 0.12 00.51±0.22	91.31 ± 0.14 00 57±0.21	91.2	7±0.10 2±0.28
1738		Datalice-scale	99.30±	0.30	99.31 ± 0.33	99.37 ± 0.31	99.0	4⊥0.06
1739		Dalik Dla al	$92.03\pm$	2.0.25	92.00±0.10	92.01 ± 0.24	91.0	4 ± 0.00
1740		Blood	/4.95±	2.95	74.85 ± 2.89	74.95±2.80	/4./	2 ± 3.04
1741		Car	$99.11\pm$	1.00	99.73 ± 0.18	99.80±0.13	99.7	8 ± 0.14
1742		Communities	85.30±	1.28	85.25 ± 0.67	85.70±0.88	85.2	5±0.89
1740		Credit-g	79.07±	5.94	/8.38±4.85	/8.66±5.36	/8.0	8±4.92
1743		Diabetes	81.86±	5.34	82.56±5.12	81.99±5.31	82.0	2 ± 4.71
1744		Eucalyptus	91.58±	0.81	91.34 ± 0.99	91.49 ± 0.88	91.7	4±0.45
1745		Heart	93.69 ±	1.26	93.45 ± 1.60	93.52 ± 1.58	93.4	2 ± 1.49
1746		Junglechess	93.29±	0.29	93.43 ± 0.31	93.64±0.36	93.4	5 ± 0.33
1747		Myocardial	$61.72 \pm$	2.65	63.64±3.08	61.95 ± 3.44	61.7	7 ± 2.43
1748		Sequence-type	$96.50 \pm$	0.92	96.44 ± 1.01	96.69±0.97	96.6	1 ± 0.91
1749		Tic-tac-toe	99.67 ±	0.36	$99.52 {\pm} 0.51$	99.59 ± 0.41	99.4	5 ± 0.56
1750		Vehicle	$96.08 \pm$	0.46	96.22±0.28	96.12 ± 0.52	96.1	5 ± 0.52
1751				(b) R	egression (RMS	E)		
1752	Dataset	M = 1	0	j	M = 20	M = 30		$M = \operatorname{All}$
1753	Bike	112.57±1	1.73	11	1.46±1.72	110.71±2.	62	111.09±1.29
1734	Crab	2.13±0.	02	2.	13±0.03	2.13±0.0	1	$2.14{\pm}0.03$
1/55	Forest-fires	83.33±23	3.21	83.	19 ± 22.47	82.17±22.	78	81.14±23.21
1756	Housing	56266.21+2	248.64	56069	9.83 ± 406.99	56047.02±12	25.19	56056.16+240.30
1757	Insurance	4622.89 ± 1	73.26	4578	$.14{\pm}149.83$	4615.54 ± 16	0.49	4614.74±196.02
1758	Solution-mi	x 0.02+0.	00	0.	02+0.01	0.02+0.0	0	0.02+0.00
1759	Wine	$0.68\pm0.$	01	0.	67±0.00	0.68 ± 0.0	Õ	0.68 ± 0.00

1782 H.7 LEARNING WITH OTHER OBJECTIVES

Table 12: Evaluation results with various loss objectives on linear model, showing (a) AUC across
15 datasets for classification and (b) RMSE across 7 datasets for regression. Best performances are
bolded.

		(a) Clas	ssification (AUC)		
Data	set	Supervised Contrastive Lea	arning CLIP	Reconstruction	Cross-entropy
Adul	t	91.32±0.12	91.29±0.12	91.36±0.13	91.26±0.13
Balar	nce-scale	99.51±0.33	99.51±0.22	99.65±0.27	99.57±0.32
Bank	2	$92.08 {\pm} 0.18$	$92.06 {\pm} 0.15$	$91.95 {\pm} 0.28$	92.19±0.23
Bloo	d	74.85 ± 2.89	74.54 ± 3.17	74.78 ± 2.56	74.96±2.98
Car	munition	99.73 ± 0.18 85.25 ± 0.67	99.74±0.17 85.53±0.42	99.79±0.18 85.22±0.64	99.77 ± 0.11 85.22±0.66
Cred	it_σ	85.25±0.07 78 38+4 85	05.55±0.42 78 20+5 74	33.33 ± 0.04 78 07+5 23	63.33 ± 0.00 78 31+5 36
Diab	etes	82.56+5.12	81.60 ± 5.45	82.33 ± 5.27	81.59 ± 5.47
Euca	lyptus	91.34±0.99	91.45±1.20	91.47±1.03	91.42±0.96
Hear	t	93.45±1.60	93.44±1.45	$93.40{\pm}1.47$	93.40±1.64
Jung	lechess	93.43±0.31	92.70 ± 0.08	93.08 ± 0.42	93.05±0.29
Myo	cardial	63.64 ±3.08	62.13 ± 3.40	60.66 ± 3.00	61.26 ± 2.60
Sequ Tic-t	ac-toe	90.44±1.01 99.52+0.51	90.05±0.77 99.49+0.32	90.41 ± 1.03 99.47 ± 0.25	90.31 ± 0.78 99.39+0.42
Vehi	cle	96.22+0.28	96.12 ± 0.52	96.15 ± 0.48	96.05±0.58
		(b) Reg	gression (RMSE)		
Dataset	Supervi	sed Contractive Learning	CUIP	Reconstruction	Cross-entropy
Bike	Supervi	111.46+1.72	112 66+2 16	112 56+2 08	112 87+2 16
Crab		2.13 ± 0.03	2.12 ± 0.02	2.12+0.03	2.12+0.02
Forest-fires		83.19 ± 22.47	87.95 ± 24.07	82.04 ± 22.03	82.45+22.59
Housing	-	56069.83±406.99	55967.38±301.06	56048.85±31.9	2 56381.03±248.59
Insurance		4578.14±149.83	4615.23±201.58	4644.91±122.8	4 4572.54±159.84
Solution-mix		$0.02{\pm}0.01$	$0.02{\pm}0.00$	$0.02{\pm}0.00$	$0.02{\pm}0.00$
Wine		$0.67{\pm}0.00$	$0.68 {\pm} 0.00$	$0.68{\pm}0.00$	$0.67{\pm}0.00$

1836 H.8 IMPACT OF THE NUMBER OF LLM TRIALS1837

Table 13: Evaluation results with various numbers of trials on linear model, showing (a) AUC across
15 datasets for classification and (b) RMSE across 7 datasets for regression. Best performances are
bolded.

			(a) Classifica	ation (AUC)			
Number	r of Trials	5	10	20	30	40	
Adult	01	22 1 1 1	01 21 + 0.00	$-$ 01 22 \pm 0.08	01.20 ± 0.08	01 22+0 12	
Ralance	91. 00	52±0.14 67±0.20	91.31 ± 0.09	9 91.32 \pm 0.08	91.29 ± 0.08 00.53 ±0.33	91.32 ± 0.12 00.51±0.33	
Bank	-scale 99.	98+0.15	92.00 ± 0.22	$8 9191\pm0.29$	99.33 ± 0.33 91.96 ±0.22	99.31 ± 0.33 92.08+0.18	
Blood	74	81+2.84	74.94 ± 2.99	75.14 ± 2.98	74.84 ± 2.90	74.85 ± 2.89	
Car	99.	56 ± 0.24	99.57±0.24	4 99.73 \pm 0.13	99.75±0.14	99.73 ± 0.18	
Commu	inities 85.	27±0.52	85.13±0.64	4 85.01±0.95	$85.18 {\pm} 0.79$	85.25±0.67	
Credit-g	g 77.	55±4.74	78.21±5.70	0 78.77±4.80	$78.60{\pm}5.35$	$78.38{\pm}4.85$	
Diabete	s 82.	$07{\pm}5.48$	81.94±5.38	8 81.76±5.16	$81.71 {\pm} 5.54$	82.56±5.12	
Eucalyr	otus 91.	12 ± 0.99	91.56 ± 0.68	8 91.71±0.71	91.32 ± 1.01	91.34 ± 0.99	
Heart	93.	69 ± 1.37	93.98±1.57	7 93.52 \pm 1.29	93.55 ± 1.27	93.45 ± 1.60	
Junglec	hess 93.4	45 ± 0.28	93.70±0.40	93.48 ± 0.32	93.60 ± 0.39	93.43 ± 0.31	
Myocar	01a1 $01.$	13 ± 3.08 50 ± 0.72	$01.0/\pm 2.00$	$\begin{array}{cccccccccccccccccccccccccccccccccccc$	04.03 ± 2.33 06.50 ± 0.71	05.04 ± 3.08 06.44 \pm 1.01	
Tic-tac-	to 0	73+018	90.30 ± 1.0	90.00 ± 1.09 90.52 ± 0.54	90.30 ± 0.71 90.34 ±0.39	90.44 ± 1.01 99.52 ± 0.51	
Vehicle	96.	17±0.44	95.96±0.3	7 96.04 \pm 0.29	96.07±0.49	96.22±0.31	
			(b) Regressi	on (RMSE)			
Number of Trials	5		10	20	30	4	.0
Bike	111.39±2.14	111	.60±2.09	111.10±1.76	112.38±2.4	3 111.46	5±1.72
Crab	$2.13 {\pm} 0.02$	2.	13 ± 0.02	$2.13{\pm}0.03$	$2.12{\pm}0.02$	2.13=	±0.03
Forest-fires	82.46±22.67	81.	84±22.78	81.56±23.50	81.75±21.9	0 83.19=	£22.47
Housing	55959.59±258.	55 56001 3 4618	1.28 ± 136.26 29+179.43	56124.43 ± 103.62 4576 85+156 01	$56134.1/\pm 310$ 4598.24 ± 186	0.96 56069.83 34 4578.14	5±406.99 +149.83
Solution-mix	0.02±0.00	0.	02 ± 0.00	0.02±0.00	0.02±0.00	0.02	±149.85
Wine	$0.68 {\pm} 0.00$	0.	$68 {\pm} 0.01$	$0.68 {\pm} 0.00$	$0.68 {\pm} 0.00$	0.67=	±0.00