

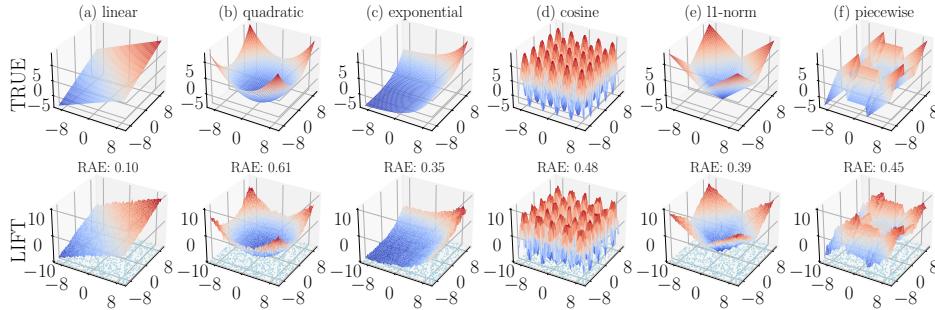
# LIFT: Language-Interfaced Fine-Tuning for Non-Language Machine Learning Tasks

Tuan Dinh\*, Yuchen Zeng\*, Ruisu Zhang, Ziqian Lin, Michael Gira, Shashank Rajput, Jy-yong Sohn, Dimitris Papailiopoulos, Kangwook Lee

University of Wisconsin-Madison, USA

## Abstract

Fine-tuning pretrained language models (LMs) without making any architectural changes has become a norm for learning various language downstream tasks. However, for *non-language* downstream tasks, a common practice is to employ task-specific designs for input, output layers, and loss functions. For instance, it is possible to fine-tune an LM into an MNIST classifier by replacing the word embedding layer with an image patch embedding layer, the word token output layer with a 10-way output layer, and the word prediction loss with a 10-way classification loss, respectively. A natural question arises: can LM fine-tuning solve non-language downstream tasks *without* changing the model architecture or loss function? To answer this, we propose **Language-Interfaced Fine-Tuning (LIFT)** and study its efficacy and limitations by conducting an extensive empirical study on a suite of non-language classification and regression tasks. LIFT does not make *any* changes to the model architecture or loss function, and it solely relies on the natural language interface, enabling “no-code machine learning with LMs.” We find that LIFT performs relatively well across a wide range of low-dimensional classification and regression tasks, matching the performances of the best baselines in many cases, especially for the classification tasks. We report the experimental results on the fundamental properties of LIFT, including its inductive bias, sample efficiency, ability to extrapolate, robustness to outliers and label noise, and generalization. We also analyze a few properties/techniques specific to LIFT, *e.g.*, context-aware learning via appropriate prompting, quantification of predictive uncertainty, and two-stage fine-tuning. Our code is available at <https://github.com/UW-Madison-Lee-Lab/LanguageInterfacedFineTuning>.



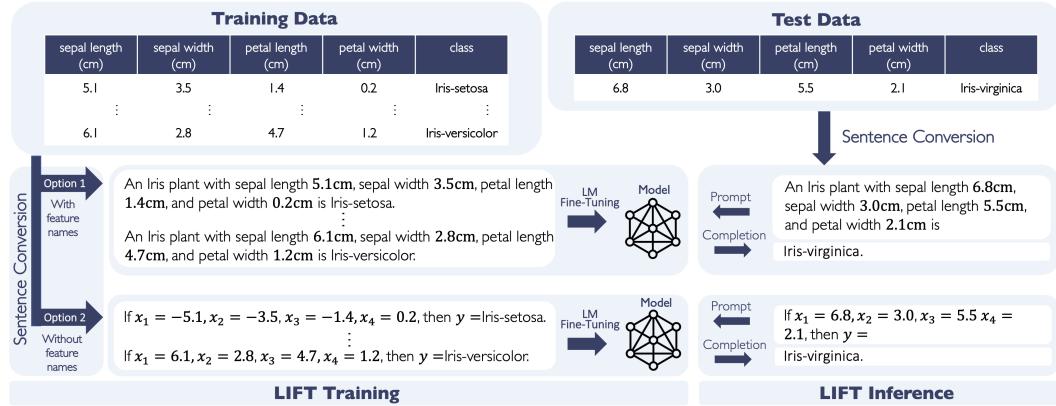
**Figure 1: Approximating various functions with Language-Interfaced Fine-Tuning (LIFT) using GPT-J.**  
We visualize the target functions (first row) and the predictor functions learned by LIFT on GPT-J (second row). Blue dots show the 1000 training samples. One can observe that LIFT well approximates the target functions.

\*Equal contribution. Emails: Tuan Dinh (tuan.dinh@wisc.edu), Yuchen Zeng (yzeng58@wisc.edu)

# Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Methodology and Experimental Setup</b>	<b>4</b>
<b>3</b>	<b>Basic Findings on LIFT</b>	<b>6</b>
3.1	Performance on Standard Machine Learning Tasks . . . . .	6
3.2	Can We Understand the Inductive Biases of Language Models via LIFT? . . . . .	7
3.3	How Many Samples Does LIFT Need? . . . . .	10
3.4	How Robust Is LIFT? . . . . .	10
3.4.1	Robustness to Outliers in Training Data . . . . .	10
3.4.2	Robustness to Label Corruption on Training Data . . . . .	12
3.4.3	Robustness to Class-Imbalance of Training Data . . . . .	12
3.4.4	Robustness to Feature Corruption on Test Data . . . . .	13
3.5	Can LIFT Interpolate and Extrapolate? . . . . .	14
3.6	Effects of Language Model Choices . . . . .	14
<b>4</b>	<b>Evaluation of LIFT-Specific Learning Properties</b>	<b>17</b>
4.1	Does LIFT Benefit from Incorporating Feature Names? . . . . .	17
4.2	Can LIFT Quantify Predictive Uncertainty via Non-deterministic Decoding? . . . . .	19
<b>5</b>	<b>LIFT Combined with Advanced Techniques</b>	<b>19</b>
5.1	Two-stage Fine-Tuning . . . . .	19
5.2	Data Augmentation . . . . .	20
<b>6</b>	<b>Related Works</b>	<b>21</b>
<b>7</b>	<b>Discussion</b>	<b>23</b>
7.1	Limitations and Open Questions . . . . .	23
7.2	Broader Impact . . . . .	24
<b>8</b>	<b>Conclusion</b>	<b>24</b>
<b>A</b>	<b>Experimental Setup</b>	<b>35</b>
A.1	Datasets . . . . .	35
A.2	LIFT and Baseline Implementation . . . . .	37
A.2.1	Pretrained Language Models and Baselines . . . . .	37
A.2.2	Computing Resources . . . . .	37
A.2.3	Hyperparameter Selection . . . . .	37
<b>B</b>	<b>Detailed and Extended Experimental Results</b>	<b>38</b>
B.1	LIFT’s Training . . . . .	38
B.2	Extended Results . . . . .	39
B.2.1	Classification and Regression Evaluations on All Baseline Methods . . . . .	39
B.2.2	Quantitative Classification Evaluations on Neural-Net-Based Synthetic Datasets	40
B.2.3	Estimating Predictive Uncertainty with LIFT/GPT-3 . . . . .	41
B.2.4	Does LIFT Benefit from Incorporating Feature Names? . . . . .	41
B.2.5	Robustness to Label Corruption . . . . .	41
B.2.6	Visualization of Regression Models . . . . .	42
B.3	Can LIFT Perform Ridge Regression via Data Augmentation? . . . . .	42

# 1 Introduction



**Figure 2: A high-level illustration of the Language-Interfaced Fine-Tuning (LIFT) framework.** LIFT has a two-phase procedure: (1) converting the dataset into sentences and (2) fine-tuning the pretrained language model (e.g., GPT) on the obtained sentences. Shown in this figure is a visualization of how LIFT can be applied to the Iris classification task. We first convert the entire Iris dataset into plain English sentences, as shown on the left. Since feature names and the task description are available for this task, one could incorporate them as part of the prompt (as option 1 in the figure). (In Sec. 4.1, we show that adding such contextual information to prompts helps LIFT achieve a higher predictive accuracy.) One may also choose to use a simpler prompt with a generic naming convention ( $x_1, x_2, \dots, x_d$ ) for  $p$  features (as option 2 in the figure). After the sentence conversion step, LIFT fine-tunes a pretrained LM with the sentence set without making any changes to model architecture or loss. At inference time, we convert the test samples to a sentence form using the same prompt, excluding the label part. LIFT performs surprisingly well in various non-language regression/classification tasks, and we summarize our main findings in Table 3. Note that to obtain a model for a given task, all we need here is to design proper sentence templates for LIFT and no changes to architecture or loss functions are needed.

Deep neural networks have been highly successful across a multitude of domains, from computer vision [1, 2] and natural language processing [3, 4], to game playing [5, 6]. Most advances in deep learning have come with a variety of domain-specific designs for network architectures, such as convolutional filters [7, 8, 9] for vision tasks, or recurrent modules [10, 11] and attention mechanisms [12, 13] in the context of natural language processing. A domain-and-modality agnostic model that can be adapted to solve tasks across different modalities and domains has become a desideratum [14], motivating great efforts in transfer learning [15] and multi-modal learning [16]. Recently, transformer-based language models (LMs) [13, 17, 18, 19] exhibited impressive versatility across different domains and modalities. They have shown great performances for various language-based tasks [20] such as question answering [21, 22], or commonsense reasoning [23]. They have also been applied to non-language modalities [18]. For instance, GPT-2 [17] pretrained on language data can be efficiently fine-tuned to perform image classification and numerical computation [18].

When downstream tasks are language-based tasks, adapting pretrained LMs can be achieved without modifying the models’ architecture. Typically, this adaptation is enabled via simple fine-tuning [24, 25, 26, 27] or in-context few-shot learning methods [28, 29]. However, not altering the architecture may pose a limitation for transferring to non-language tasks. As their input and output formats are not in some language form, adapting LMs to these domains may seem to require architectural changes. Indeed, it has been a common practice to re-design the input/output layers and loss function to accommodate a different predictive task. For instance, to adapt GPT-2 [21] to other modalities, the frozen pretrained transformer [18] adds new input/output layers to handle different types of input/output. To make such changes, one must have a good understanding of the underlying principles of LMs and an ability to make proper modifications at the code level.

A natural question that arises is whether such changes are necessary. In other words,

Does language model fine-tuning work for non-language tasks  
**without** changing the architecture or loss function at all?

To answer this question, we consider a simple fine-tuning procedure for pretrained LMs, which we refer to as **Language-Interfaced Fine-Tuning (LIFT)**. The LIFT framework can be used to learn a predictor for any classification or regression task. LIFT runs in two phases: (1) it converts labeled samples into sentences, and (2) it fine-tunes pretrained LMs on the sentence dataset without altering the architecture or loss function.

Fig. 2 illustrates how we fine-tune GPT with LIFT to solve the Iris classification task [30]. LIFT first converts each labeled sample into a sentence. There are two options here. The first option is to incorporate the feature names and the task description into the sentence template. In this example, we could convert a training sample  $\mathbf{r}$  into “An Iris plant with sepal length  $\mathbf{r}.\text{sepal\_length}$ , sepal width  $\mathbf{r}.\text{sepal\_width}$ , petal length  $\mathbf{r}.\text{petal\_length}$ , and petal width  $\mathbf{r}.\text{petal\_width}$  is  $\mathbf{r}.\text{class}$ .” Here, we use the dot notation, *i.e.*,  $\mathbf{r}.\star$  denotes the string conversion of the corresponding attribute of sample  $\mathbf{r}$ . One may also adopt a simpler (and more generic) sentence template such as “If  $x_1=\mathbf{r}.\mathbf{x1}$ ,  $x_2=\mathbf{r}.\mathbf{x2}$ , ...,  $x_p=\mathbf{r}.\mathbf{xp}$ , then  $y=\mathbf{r}.\mathbf{y}$ ,” if there are  $p$  features. We then fine-tune LMs without changing either architecture or loss function. Once the model is fine-tuned, we perform inference as follows. LIFT first converts test samples into sentences using the same template while leaving the part that needs to be predicted empty. It then feeds the converted sentences as prompts to the fine-tuned model. The output tokens are parsed to provide the final predictions.

Our work empirically shows that LIFT can provide high accuracy solutions for a variety of non-language tasks. In Fig. 1, we show examples of real functions learned by GPT-J models [31] fine-tuned using LIFT, when 1000 training samples are given. Recall that LIFT does not require any changes in the architecture or loss function. Thus, our findings show that such changes to architecture/loss function might *not* be necessary, even when the target predictive task is not a language task. Thus, LIFT can be almost perceived as a “no-code machine learning” framework as the data-to-sentence conversion is extremely straightforward even without extensive programming skills and machine learning backgrounds.

Motivated by these intriguing properties, we investigate the efficacy and limitations of LIFT on non-language tasks by conducting an extensive empirical study on a suite of classification and regression tasks. *First*, we observe that LIFT performs well across a wide range of low-dimensional classification and regression tasks. In most cases, it nearly matches (or slightly outperforms) the performance of the best baseline models. To further understand LIFT, we conduct experiments testing the fundamental learning properties, *e.g.*, its inductive bias, sample efficiency, ability to extrapolate, and worst- and average-case noise robustness. *Third*, we also study a few unique properties specific to LIFT, *e.g.*, context-aware learning with task-specific prompting and quantifying the predictive uncertainty via non-deterministic decoding. *Lastly*, to improve upon the basic fine-tuning, we employ a few different techniques: a two-stage fine-tuning with synthetic pretext tasks and then data augmentation with mixup. We find that the two-stage fine-tuning improves the performance of LIFT in the regime of low training samples. We discuss the current limitations of LIFT and future investigations towards making LIFT more effective and efficient.

**Scope of study** The goal of this work is to understand: (i) what LIFT can and cannot do, (ii) properties of fine-tuned models obtained via LIFT, and (iii) the behavior of LIFT when combined with advanced techniques. We emphasize that our goal is *not* to achieve the state-of-the-art performance via LIFT but to provide a thorough investigation of this new approach.

## 2 Methodology and Experimental Setup

We describe the methodology used in our empirical study.

**LIFT training.** To fine-tune a pretrained LM with LIFT on a target supervised learning task, we apply two steps: (1) convert each sample into a sentence with a fixed template, and (2) fine-tune LMs with sentence datasets. We use the default cross-entropy loss for token prediction in LMs. Our generic template (without feature names and task description) for sample  $\mathbf{r}$  is

When we have  $x_1=\mathbf{r}.\mathbf{x1}$ ,  $x_2=\mathbf{r}.\mathbf{x2}$ , ...,  $x_p=\mathbf{r}.\mathbf{xp}$ , what should be  $y$ ?  $\underbrace{\quad\quad\quad}_{\text{question}}$   $\underbrace{\#\#\#}_{\text{q/a separator}}$   $\underbrace{y = \mathbf{r}.\mathbf{y}}_{\text{answer}}$   $\underbrace{@@@@}_{\text{end of answer}}$ ,

if  $\mathbf{r}$  has  $p$  attributes. Here, we use the separator convention recommended by OpenAI [32] – “ $\#\#\#$ ” for question/answer separation, and “ $@@@@$ ” for end of generation. When attributes names and task

descriptions are available, one can use a more informative prompt, as we showed in Fig. 2. For all such cases, we provide the actual prompts in Sec. 4.1. We report the learning curves of LIFT on several classification and regression tasks in Appendix B.1.

**LIFT inference.** For inference, we use the same prompt template except for the answer and end-of-answer parts. Once the fine-tuned LM completes the test prompt, we simply parse the output tokens. For the case of classification, we simply compare the generated text with the string representation of the class names. For the case of regression, we convert the generated string into a number. For instance, with the output being “ $y=10.35@\text{extratokens}$ ”, we split the output sentence by the stop separator “ $@@@$ ” into two parts. Taking the first part “ $y=10.35$ ”, we parse the value “10.35” as our final prediction.

The generated output might be *invalid*. For the case of classification, the output string may not match any of the actual classes. If this happens, we simply declare a misclassification. Note that one could always return the closest class using a string metric or word similarity to obtain slightly better accuracy. For the case of regression, the generated output is considered invalid if the string-to-number parsing fails. When this happens, we increase the decoding temperature [33, 34, 35], which controls the randomness in the generation. In particular, we change the generation mode from deterministic (temperature 0) to random (temperature 0.75) and then repeat the inference up to five times. If all attempts fail, then we simply return the average value of the training set. Note that invalid output occurs very rarely and is less than or around 1% in most tested cases.

As for evaluation metrics, we measure accuracy for classification tasks, and RMSE, RAE errors for regression tasks, where  $\text{RAE} = \sum_{i=1}^n |\hat{y}_i - y_i| / \sum_{i=1}^n |\frac{1}{n} \sum_{j=1}^n y_j - y_i|$  and  $\text{RMSE} = \sqrt{\frac{\sum_{i=1}^n (\hat{y}_i - y_i)^2}{n}}$  on each dataset  $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$  with  $n$  samples, features  $\mathbf{x} \in \mathcal{X} \subset \mathbb{R}^p$ , and outcome  $y$ .

**Pretrained LMs.** We apply LIFT on two pretrained LMs: GPT-J [31] and GPT-3 [19]. For GPT-J, we apply LoRA [24] for parameter-efficient fine-tuning, which constrains the weight matrix updates to be low-rank. For experiments on GPT-J, we used p3.8xlarge and p3.2xlarge instances from AWS and RTX3090 GPUs in the local server. Since GPT-3 is not fully publicly available, we use the API provided by OpenAI to perform black-box GPT-3 finetuning. More details are in Appendix A.2.1.

**Datasets.** We design and select a wide range of datasets to better understand the behavior of LIFT. For *classification*, we use three types of non-language data: low-dimensional synthetic datasets, real tabular datasets in OpenML [36], and vision datasets (MNIST [37], Fashion-MNIST [38] and their permuted variants [39]). For *regression*, we use both synthetic and real datasets. For synthetic datasets, we defined samples  $(\mathbf{x}_i, y_i)$  of input-output pair as  $y \sim f(\mathbf{x}) + \mathcal{N}(0, \sigma^2)$ , where  $\sigma^2 \geq 0$  is the noise level. Unless otherwise stated, we sample the feature  $\mathbf{x}$  uniformly from a hypercube  $[L, U]^p$ , where  $L$  and  $U$  are minimum/maximum feature values, and  $p$  is the number of features. Following the suggestion by [40], we consider various functions  $f$  for regression tasks: (i) linear function, (ii) quadratic function, (iii) exponential function, (iv) cosine function, (v)  $\ell_1$ -norm function, and (vi) piece-wise linear function. Their 2D visualizations are provided in the first row of Fig. 1. We also use four real datasets: Medical Insurance (Insurance) [41], Combined Cycle Power Plant (CCPP) [42], Servo [43], and Student Performance (Student) [44]. More details are included in Appendix A.1.

**Baselines.** We consider standard machine learning algorithms [45, 46]. The classification baselines are logistic regression (*LogReg*), decision tree (*DT*), k-nearest neighbor (*KNN*), support vector machine with Gaussian kernel (*SVM*), a four-layer ReLU neural network (*MLP*) with 200 neurons for each hidden layer, random forest (*RF*), and XGBoost (*XG*). We also compared an optimal deterministic classifier (*ODC*) that outputs the most dominant class. For regression, we use polynomial regression (*PR*), kernel ridge regression (*KR*) with radial basis function kernel,  $k$ -nearest neighbors (*KNN*), a three-layer ReLU neural network (*MLP*) with 50 hidden neurons at each layer, Gradient Boosting Trees (*GBT*), random forest (*RF*), and Gaussian process (*GP*). To select hyperparameters for models, we apply the grid search on a set of parameters’ values and use cross-validation on the training set to choose the best configurations (see details in Appendix A.2.1).

Topic	Findings
Overall performance	On various classification tasks, LIFT achieves accuracies comparable to strong baselines (see Table 4). For regression tasks, LIFT well approximates different types of low-dimensional functions (see Fig. 1), while LIFT does not perform well for high-dimensional cases (see Table 28).
Robustness	For regression tasks, LIFT is robust to outliers in training data (Fig. 10). For classification tasks, LIFT has the performance comparable to baselines when training data has corrupted labels (Fig. 11) or class imbalance (Table 12). LIFT is vulnerable to adversarial examples (test-time feature corruption) transferred from simple neural networks (Table. 13).
Context-aware learning	Performance of LIFT on classification tasks can be improved by designing the prompt in a way that the feature names and the target task are specified. The improvement is significant when the description of the feature names and the target task can be interpreted with common knowledge (Table 18).
Two-stage training	Warming up LIFT with pretext tasks using synthetic data improves the prediction performance, especially in the low-data regime (Fig. 20).
Data augmentation	For classification tasks, training with augmented data significantly improves the tolerance of LIFT against perturbed test data (Table 21).

Table 3: Summary of the main findings.

### 3 Basic Findings on LIFT

In this section, we empirically study LIFT on standard non-language classification and regression tasks. Table 3 summarizes our main findings on LIFT in terms of classification/regression performance, robustness of LIFT against outlier/corruption/attack, as well as the effect of context-aware learning (Sec. 4), two-stage fine-tuning, and data augmentation on LIFT (Sec. 5).

#### 3.1 Performance on Standard Machine Learning Tasks

**Classification.** Consider a  $c$ -way classification problem. Table 4 compares classification accuracies of evaluated algorithms on a wide range of tasks. Overall, we observe that LIFT can perform comparably well across most datasets compared to the best-performing baseline algorithms (XGBoost, RBF-SVM, and MLP). In most cases, LIFT ranks highly in the top three best-performing methods. For instance, LIFTs achieve the best classification accuracies on the synthetic dataset 9Clusters (100%), the real dataset OpenML-28 (98.99%), and the vision dataset MNIST (98.15%). We provide the full comparison of LIFT with all other baselines in Table 27 (Appendix B.2.1).

On the synthetic datasets, LIFT achieves the highest accuracies in many cases and performs competitively compared to the best baselines across all the tasks. On the tabular OpenML data, the performances of LIFT are comparable to the strongest baselines, such as XGBoost, RF, and RBF-SVM, across the datasets. As the difficulty of tasks varies, which can be estimated by the average performance of baselines, LIFT also suffers from the performance degradation. LIFT can perform relatively well even when the number of features is as large as hundreds, though the number of features LIFT can input is restricted due to the limited input context length of LMs. However, when the number of classes is large, both LIFT/GPT-J and LIFT/GPT-3 suffer from lower accuracies than many baselines, even though they manage to be better than ODC. For instance, on the Margin dataset (OpenML ID=1491) with 100 classes, the accuracy difference between LIFT/GPT-3 and the best algorithm (RBF-SVM) is nearly 20%.

More interestingly, LIFT achieves relatively high accuracies on vision datasets, and these results are comparable to the performances of several popular convolutional neural networks (CNNs) [2, 47]. In particular, LIFT/GPT-3 achieves 98.15% and 90.18% on MNIST and Fashion MNIST. For some comparisons, recent work [48] reports the accuracies on MNIST and Fashion MNIST as 98.81% and 86.43% for AlexNet [49], 99.37% and 94.39% for ResNet [8], 99.57% and 90.03% for CapsuleNet [50]. For the permuted MNIST and permuted Fashion MNIST, LIFT/GPT-3 obtains

Dataset (ID)	$p / c$	ODC	LogReg	DT	RBF-SVM	XG		LIFT/GPT-J	LIFT/GPT-3
<b>Synthetic Data</b>									
circles (3)	2 / 2	50.00	48.58 $\pm$ 1.94	77.42 $\pm$ 0.24	<b>83.08<math>\pm</math>0.59</b>	81.42 $\pm$ 0.31		79.95 $\pm$ 1.53	81.17 $\pm$ 0.42
two circles (6)	2 / 2	50.00	49.83 $\pm$ 4.18	75.50 $\pm$ 0.20	80.00 $\pm$ 0.54	79.25 $\pm$ 0.35		75.92 $\pm$ 1.65	<b>81.42<math>\pm</math>0.82</b>
blobs (2)	2 / 4	25.00	<b>96.75<math>\pm</math>0.00</b>	96.08 $\pm$ 0.82	<b>96.75<math>\pm</math>0.00</b>	96.17 $\pm$ 0.12		96.17 $\pm$ 0.59	96.67 $\pm$ 0.24
moons (4)	2 / 4	50.00	88.58 $\pm$ 0.12	99.25 $\pm$ 0.41	<b>100.00<math>\pm</math>0.00</b>	99.83 $\pm$ 0.12		99.58 $\pm$ 0.42	<b>100.00<math>\pm</math>0.00</b>
9Clusters (1)	2 / 9	11.25	100.00 $\pm$ 0.00	100.00 $\pm$ 0.00	<b>100.00<math>\pm</math>0.00</b>	<b>100.00<math>\pm</math>0.00</b>		99.75 $\pm$ 0.00	<b>100.00<math>\pm</math>0.00</b>
<b>Tabular Data (OpenML)</b>									
Customers (1511)	8 / 2	68.18	<b>87.12<math>\pm</math>0.54</b>	85.98 $\pm$ 0.53	86.36 $\pm$ 0.00	85.23 $\pm$ 0.00		85.23 $\pm$ 1.61	84.85 $\pm$ 1.42
Pollution (882)	15 / 2	50.00	58.33 $\pm$ 11.79	<b>77.78<math>\pm</math>3.93</b>	58.33 $\pm$ 6.81	63.89 $\pm$ 7.86		63.89 $\pm$ 3.93	63.89 $\pm$ 7.86
Spambase (44)	57 / 2	60.59	93.27 $\pm$ 0.00	90.7 $\pm$ 0.14	93.70 $\pm$ 0.00	<b>95.87<math>\pm</math>0.00</b>		94.03 $\pm$ 0.54	94.90 $\pm$ 0.36
Hill-Valley (1479)	100 / 2	49.79	77.78 $\pm$ 0.00	56.38 $\pm$ 0.89	68.72 $\pm$ 0.00	59.26 $\pm$ 0.00		<b>100.00<math>\pm</math>0.20</b>	<b>99.73<math>\pm</math>0.19</b>
IRIS (61)	4 / 3	33.33	96.67 $\pm$ 0.00	97.77 $\pm$ 3.85	<b>100.00<math>\pm</math>0.00</b>	<b>100.00<math>\pm</math>0.00</b>		96.67 $\pm$ 0.00	97.0 $\pm$ 0.00
TAE (48)	5 / 3	35.48	45.16 $\pm$ 4.56	65.59 $\pm$ 5.49	53.76 $\pm$ 6.63	<b>66.67<math>\pm</math>8.05</b>		61.29 $\pm$ 6.97	65.59 $\pm$ 6.63
CMC (23)	9 / 3	42.71	49.49 $\pm$ 0.83	56.72 $\pm$ 0.32	56.50 $\pm$ 0.97	52.43 $\pm$ 0.42		49.83 $\pm$ 0.28	<b>57.74<math>\pm</math>0.89</b>
Wine (187)	13 / 3	38.89	<b>100.00<math>\pm</math>0.00</b>	93.52 $\pm$ 2.62	<b>100.00<math>\pm</math>0.00</b>	97.22 $\pm$ 0.00		93.52 $\pm$ 1.31	92.59 $\pm$ 1.31
Vehicle (54)	18 / 4	25.88	80.39 $\pm$ 1.00	63.92 $\pm$ 2.37	<b>81.18<math>\pm</math>0.48</b>	73.14 $\pm$ 0.28		64.31 $\pm$ 2.37	70.20 $\pm$ 2.73
LED (40496)	7 / 10	11.00	68.67 $\pm$ 0.94	66.33 $\pm$ 2.87	68.00 $\pm$ 0.82	66.00 $\pm$ 0.82		65.33 $\pm$ 0.47	<b>69.33<math>\pm</math>2.05</b>
OPT (28)	64 / 10	10.14	96.53 $\pm$ 0.22	89.8 $\pm$ 1.09	97.95 $\pm$ 0.00	97.48 $\pm$ 0.17		98.22 $\pm$ 0.11	<b>98.99<math>\pm</math>0.30</b>
Mfeat (12)	216 / 10	10.00	97.67 $\pm$ 0.12	87.67 $\pm$ 1.05	<b>98.83<math>\pm</math>0.24</b>	96.75 $\pm$ 0.00		94.17 $\pm$ 1.75	93.08 $\pm$ 0.24
Margin (1491)	64 / 100	0.94	81.35 $\pm$ 0.15	43.86 $\pm$ 1.21	<b>81.98<math>\pm</math>0.30</b>	70.21 $\pm$ 0.29		50.23 $\pm$ 1.33	59.37 $\pm$ 0.92
Texture (1493)	64 / 100	0.94	81.67 $\pm$ 0.97	46.88 $\pm$ 1.93	<b>83.44<math>\pm</math>0.89</b>	70.73 $\pm$ 1.41		50.32 $\pm$ 2.18	67.50 $\pm$ 1.42
<b>Image Data</b>									
MNIST		11.35	91.95 $\pm$ 0.69	87.42 $\pm$ 0.64	97.70 $\pm$ 0.97	97.69 $\pm$ 0.04		97.01 $\pm$ 1.15	<b>98.15<math>\pm</math>0.67</b>
Permuted MNIST	784 / 10	11.35	92.58 $\pm$ 0.04	87.87 $\pm$ 0.69	<b>98.06<math>\pm</math>0.31</b>	97.62 $\pm$ 0.09		95.80 $\pm$ 0.07	96.25 $\pm$ 0.35
Fashion MNIST		10.00	85.59 $\pm$ 0.09	80.52 $\pm$ 0.40	<b>90.59<math>\pm</math>0.02</b>	90.19 $\pm$ 0.04		85.10 $\pm$ 0.19	90.18 $\pm$ 0.12
Permuted F-MNIST		10.00	84.95 $\pm$ 0.84	79.91 $\pm$ 0.93	88.04 $\pm$ 1.69	<b>89.93<math>\pm</math>0.14</b>		82.25 $\pm$ 0.27	88.92 $\pm$ 0.71

**Table 4: Accuracies ( $\uparrow$ ) on various classification datasets.** We provide the full comparison with all other baselines (KNN, MLP, and Random Forest) in Table 27 (Appendix B.2.1). We evaluate LIFT/GPTs on different classification datasets: 2D synthetic datasets, tabular datasets in OpenML [36], and image datasets, varying number of features ( $p$ ) and number of classes ( $c$ ). Overall, LIFT/GPTs perform relatively well across tasks. LIFT/GPTs can be adapted well to non-linear datasets (circles, two circles), beyond the capacity of logistic regression. On the OpenML data, LIFT/GPTs achieve competitive performances with the best methods, such as XGBoost or RBF-SVM. The performance of LIFT degrades as the number of classes is large, e.g., when the number of classes  $c=100$ . On the vision data, LIFT/GPTs perform relatively well, achieving highly competitive accuracies on both MNIST and Fashion MNIST. We note that the classes of MNIST are not fully balanced, thus ODC achieves 11.35% instead of 10% as ODC returns the optimal class learned from the training dataset.

slightly lower accuracies than the ones on the original datasets, with 96.25% and 88.92%, respectively. These results are comparable to the best-reported performances of MLPs [51], which are 98.6% and 90.9% on the two permuted variants.

We also note that LIFT can learn non-linear decision boundary: LIFT achieves 80.17% accuracy on the circles dataset while logistic regression (linear classifier) fails to perform better than ODC (50%). Moreover, while most baseline algorithms require scaling and normalization of features for achieving good performances, LIFT can directly use raw values.

**Regression.** Shown in Fig. 1 are the predictor functions learned by LIFT with 1000 samples. One can also compare the predictors learned with different models in Fig. 37. See Fig. 36 for similar results when 200 training samples are used.

Table 5 reports the quantitative performance of LIFT and other baselines on synthetic datasets with 200 training samples under low and high-dimensional cases, respectively. The full results are deferred to Table 28 in Sec. B.2.1 of Appendix. For the low-dimensional cases, the performances of LIFT are comparable to most baselines but still worse than the strongest baselines, such as GP. For high-dimensional cases, LIFT and the baselines fail to perform well on approximating all functions. This is as expected since GPT fine-tuning is based on the classification loss, which fails to capture the magnitude of the prediction error. We provide more discussion on the difficulties of regression tasks in Sec. 7.1.

### 3.2 Can We Understand the Inductive Biases of Language Models via LIFT?

In this section, we investigate the inductive biases of pretrained LMs via analyzing their decision boundaries learned with LIFT on several classification tasks. We also compare, both qualitatively

Dataset \ Method		KNN	MLP	GBT	GP	LIFT/GPT-J	LIFT/GPT-3
Linear	$p = 1$	$0.04 \pm 0.0$	$0.03 \pm 0.0$	$0.05 \pm 0.0$	$0.01 \pm 0.0$	$0.08 \pm 0.0$	$0.06 \pm 0.0$
	$p = 2$	$0.12 \pm 0.0$	$0.04 \pm 0.0$	$0.12 \pm 0.0$	$0.01 \pm 0.0$	$0.12 \pm 0.0$	$0.19 \pm 0.0$
Quadratic	$p = 1$	$0.05 \pm 0.0$	$0.03 \pm 0.0$	$0.06 \pm 0.0$	$0.01 \pm 0.0$	$0.11 \pm 0.0$	$0.13 \pm 0.0$
	$p = 2$	$0.17 \pm 0.0$	$0.06 \pm 0.0$	$0.15 \pm 0.0$	$0.02 \pm 0.0$	$0.28 \pm 0.1$	$0.22 \pm 0.0$
Exponential	$p = 1$	$0.05 \pm 0.0$	$0.02 \pm 0.0$	$0.05 \pm 0.0$	$0.01 \pm 0.0$	$0.11 \pm 0.0$	$0.09 \pm 0.0$
	$p = 2$	$0.13 \pm 0.0$	$0.07 \pm 0.0$	$0.09 \pm 0.0$	$0.04 \pm 0.0$	$0.19 \pm 0.0$	$0.20 \pm 0.0$
Cosine	$p = 1$	$0.14 \pm 0.0$	$0.38 \pm 0.1$	$0.15 \pm 0.0$	$0.04 \pm 0.0$	$0.38 \pm 0.1$	$0.44 \pm 0.1$
	$p = 2$	$0.83 \pm 0.1$	$1.06 \pm 0.0$	$0.41 \pm 0.0$	$0.31 \pm 0.0$	$0.82 \pm 0.2$	$0.65 \pm 0.1$
L1norm	$p = 1$	$0.05 \pm 0.0$	$0.03 \pm 0.0$	$0.06 \pm 0.0$	$0.03 \pm 0.0$	$0.10 \pm 0.0$	$0.09 \pm 0.0$
	$p = 2$	$0.19 \pm 0.0$	$0.06 \pm 0.0$	$0.15 \pm 0.0$	$0.07 \pm 0.0$	$0.24 \pm 0.0$	$0.20 \pm 0.0$
Piecewise	$p = 1$	$0.08 \pm 0.0$	$0.08 \pm 0.0$	$0.06 \pm 0.0$	$0.10 \pm 0.0$	$0.15 \pm 0.0$	$0.17 \pm 0.0$
	$p = 2$	$0.33 \pm 0.0$	$0.20 \pm 0.0$	$0.19 \pm 0.0$	$0.29 \pm 0.0$	$0.40 \pm 0.1$	$0.40 \pm 0.1$

**Table 5: Comparison of regression methods and LIFT in approximating various functions under low-dimensional cases.** The regression performance is measured by RAE ( $\downarrow$ ), and  $p$  is the number of features. LIFT can approximate different types of functions in low-dimensional cases ( $p = 1, 2$ ), although it fails to achieve performance comparable to that of strong baselines.

and quantitatively, the decision boundary learned by LIFT with those of baseline algorithms to find whether LIFT shares similar patterns to any of these baselines.

**Visualizing decision boundary** We construct datasets with various classification complexities to investigate the adaptability of LIFT. In particular, we construct three binary classification datasets, a 3-class dataset and a 5-class dataset (shown in the first column of Fig. 6a, Fig. 6b, and Fig. 6c). We call these datasets *neural-net-based synthetic datasets* since we generate them using a 2-layer neural network. See Fig. 22 and Appendix A.1 for detailed explanations of how we generated these datasets.

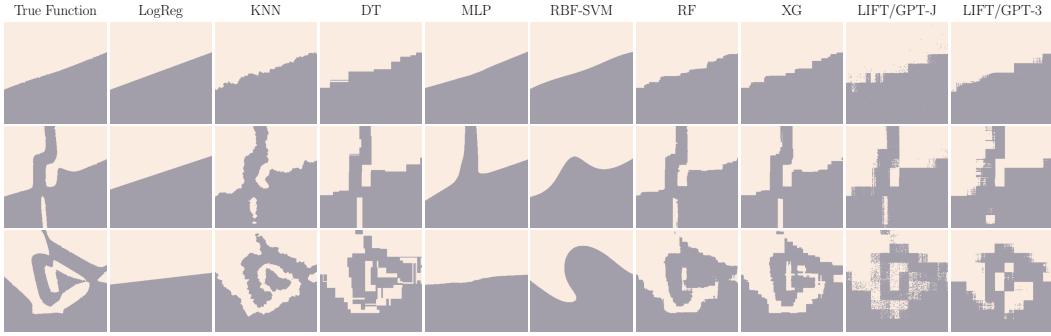
Fig. 6 visualizes the decision boundaries of models trained on the neural-net-based synthetic data. In addition, we also visualize the decision boundaries of models trained on the *label-corrupted* versions of three binary classification datasets, with the corruption probabilities being 5% and 20% (see details in Sec. 3.4.2), shown in Fig. 6d and Fig. 6e. Specifically, we consider the binary classification tasks and flip the training data labels with the provided probabilities. Overall, both GPT-J and GPT-3 models fine-tuned with LIFT can adapt well to different boundaries. They can capture the rough shapes of the decision boundaries in all three settings.

Our results indicate that fine-tuning GPT using LIFT may extract valuable features from labeled datasets to adapt to different decision boundaries, which is partly consistent with recent findings on inductive biases of language models [52].

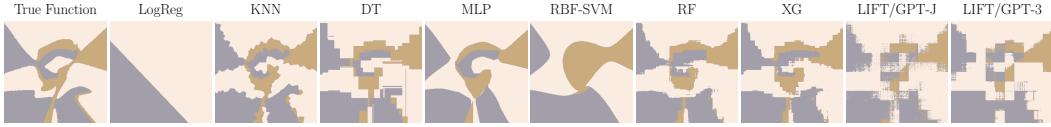
When the level of corruption increases, decision trees and XGboost are the most affected baselines. While roughly capturing the boundary, LIFT/GPT-J also shows more noisy predictions. In contrast, LIFT/GPT-3 displays great robustness against the corrupted labels while capturing the correct boundary shapes. Nevertheless, this experiment indicates the different behaviors of LIFTs from the baseline algorithms and their adaptability to different types of the decision boundary.

One interesting observation here is that LIFT’s decision boundaries are axis-parallel and show a lot of fractals. The axis-parallel boundary looks similar to the boundary of tree-based classifiers, and the fractal shapes of LIFT’s boundaries are similar to the observations on the decision boundaries of some convolution neural networks [53]. However, the main reason why LIFT’s decision boundary has such patterns seems to be due to the way it interprets numbers. Since we rely solely on the language interface, there are some artifacts due to the decimal numeral system. For instance, 0.98 and 0.99 are only one-character different, but 0.99 and 1.00 are three-characters different. We believe that such an artifact is the reason behind axis-parallel decision boundaries and fractal-like patterns.

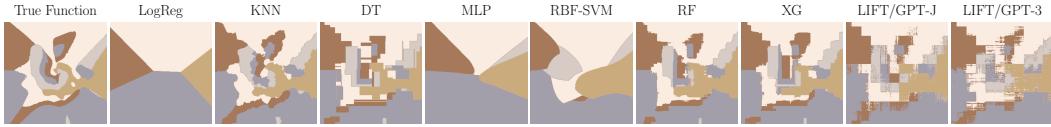
**Quantifying the similarity of decision boundaries.** To further verify whether LIFTs behave similarly to any standard algorithm, we quantify the similarity between the decision boundaries of LIFT/GPT-3 and those of the baselines. Specifically, the similarity score is the percentage of the exact classification matches between LIFT/GPT-3 and the compared method. We randomly sample



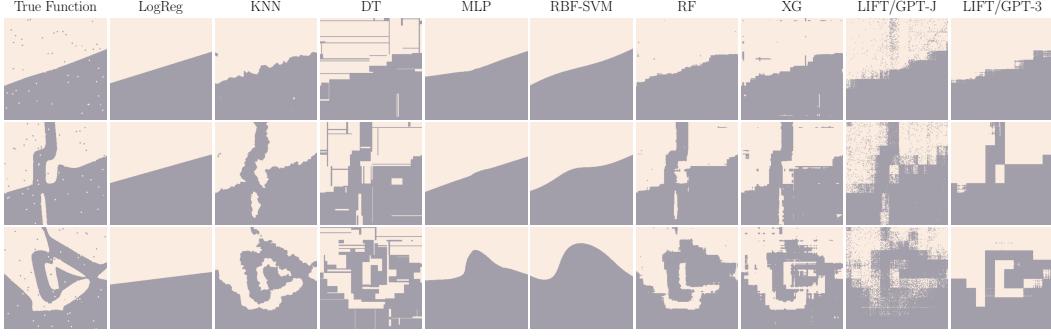
**(a) Binary classification**



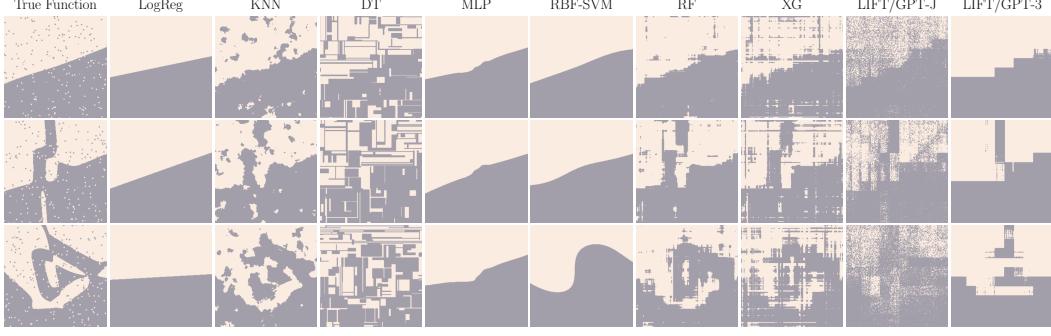
**(b) 3-way classification**



**(c) 5-way classification**



**(d) Binary classification with label corruption (corruption ratio is 5%)**



**(e) Binary classification with label corruption (corruption ratio is 20%)**

**Figure 6: Classification decision boundary visualizations on neural-net-based synthetic datasets.** The first column of each row shows the decision boundary of the training datasets. In (a), (b), and (c), we visualize the decision boundaries of models trained on datasets with two, three, and five classes. In (d) and (e), we consider the label-corrupted version of the binary-class datasets, with the corruption probabilities of 5% and 20%. We find that LIFT/GPTs adapt well and roughly estimate the true decision boundaries. The shapes of LIFT's decision boundary are likely to be axis-parallel and show multiple fractals.

Dataset (ID)	Similarity \ Method	SVM (kernel)			LogReg			KNN ( $k$ )			DT (depth $D$ )		MLP (width $W$ )			XG	RF (# estimators $E$ )		
		poly	rbf	sigmoid		K=1	K=3	K=5	D=3	D=5	W=10	W=100	W=200		E=20	E=50	E=100		
9clusters (1)		100.00	100.00	100.00	100.00	100.00	100.00	100.00	76.00	100.00	100.00	100.00	100.00	97.50	100.00	100.00	100.00		
blobs (2)		98.50	97.50	92.00	97.50	94.00	95.50	96.00	94.50	91.00	97.00	97.00	97.00	93.50	94.50	94.00	94.00		
circles (3)		54.00	93.50	48.00	51.00	85.50	88.50	88.50	67.00	80.00	84.00	92.50	92.50	87.50	85.00	87.50	89.00		
moons (4)		90.50	97.50	74.50	87.50	99.00	99.00	99.00	91.00	98.50	92.50	98.50	98.50	97.50	95.50	94.50	96.00		
two circles (6)		63.00	62.50	48.50	62.50	59.50	58.00	59.00	57.50	59.50	60.50	64.00	65.00	58.50	64.00	63.50	65.00		
CMC (23)		59.50	61.00	68.00	63.50	50.50	52.50	51.00	65.00	72.00	62.50	53.50	53.50	63.00	56.00	55.50	59.00		
Pollen (871)		66.00	74.50	66.00	67.50	60.00	63.00	64.50	69.50	66.00	60.00	59.50	58.00	63.50	64.50	62.50	67.00		
Climate (1467)		100.00	100.00	100.00	94.50	92.50	98.50	100.00	94.50	91.00	98.00	97.50	97.50	94.50	99.50	99.50	100.00		
LED (40496)		88.00	87.00	82.50	91.50	74.00	76.50	84.50	69.00	84.50	85.50	94.50	91.50	92.00	91.50	91.50	90.50		
Average		79.94	<b>85.94</b>	75.5	79.5	79.4	81.28	82.50	76.00	82.50	82.22	<b>84.11</b>	83.72	[83.06]	83.39	81.17	<b>84.5</b>		

**Table 7: Quantifying the similarity between decision boundaries of LIFT/GPT-3 and those of various baselines.** We use different settings of the baselines, where their hyperparameters are given with the baseline name, and the selected values of hyperparameters are specified in the second line. Each column reports the matching accuracy ( $\uparrow$ ) between the predictions of LIFT/GPT-3 with those of the baseline. Each score is a percentage similarity of both LIFT/GPT-3 and the baseline classifying a point with the same class. For example, a score of 100 for model A signifies that LIFT/GPT-3 classified all sampled test points in the same manner as model A, regardless of their true dataset accuracy. The last row reports the average matching accuracy. We highlight the highly matched algorithms, namely RBF-SVM, MLP ( $W=100$ ), and Random Forest ( $E=100$ ).

two sets of 200 data points from the original dataset for training and evaluation, respectively, and the results of all methods are reported in Table 7. Based on the similarity score, while we observe no similar discernible pattern between LIFT/GPT-3 and the baselines, we find that LIFT/GPT-3 appear to share the most similar behavior pattern to RBF-SVM, random forest ( $E=100$ ), and MLP ( $W=100$ ).

### 3.3 How Many Samples Does LIFT Need?

We investigate whether LIFT is sample efficient. Fig. 8 shows the sample complexity evaluation on classification and regression tasks. We find that the GPT model can be quickly fine-tuned to new tasks when LIFT is applied. For *classification*, when the number of classes increases (from the left to the right column in Fig. 8a), LIFT does need more samples for adaptation. This is probably because LIFT requires more samples to learn more complex input and output spaces of the data. For *regression*, we find that  $n = 1000$  samples are sufficient for LIFT to have small RMSE, similar to other baselines. There exist some functions (quadratic, cosine, piecewise) where LIFT has lower sample complexity than popular baselines.

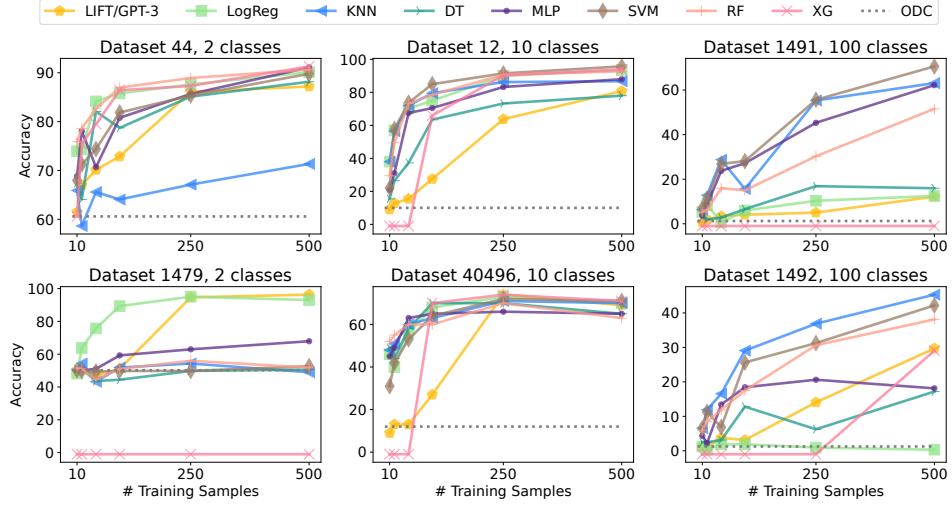
**LIFT versus in-context learning.** The in-context few-shot learning methods aim to infer new tasks by conditioning on a few training examples without fine-tuning [54, 19]. Since the context length of LMs is limited, for a new target task requiring extensive training data, fine-tuning LMs is a more appropriate choice than in-context learning. However, it is unclear if this comparison still holds when we are given only a small number of training samples. We investigate this question by comparing LIFT with the in-context learning approach on the same pretrained LMs across various data. Table 9 shows the comparison of six datasets and two GPT models. First, we note that LIFTs trained on full data consistently achieve the best performances overall. When the same number of samples is used, using LIFT is a better choice than the in-context learning methods in most cases for GPT-J, while LIFT and the in-context learning method have comparable performances for GPT-3. Note that the number of prompts for each dataset is selected to fit in the maximum context length that OpenAI GPT-3 can accept.

### 3.4 How Robust Is LIFT?

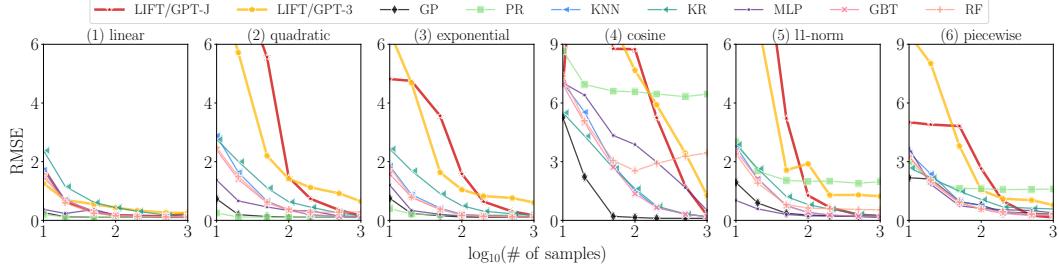
We investigate to which extent LIFT is robust to practical settings of training data having outlier samples, corruption on labels, and class imbalance. We also explore whether LIFT models are robust against test-time data corruptions, *e.g.*, random noise or adversarial perturbation [55].

#### 3.4.1 Robustness to Outliers in Training Data

We investigate the robustness of LIFT to outliers in regression tasks whose outcome  $y$  is not consistent in terms of fitting ( $x, y$ ). Fig. 10a compares the RAE values of LIFT and baselines with and without



**(a)** Classification tasks on OpenML tabular datasets. The classification performance is measured in terms of accuracy ( $\uparrow$ ).

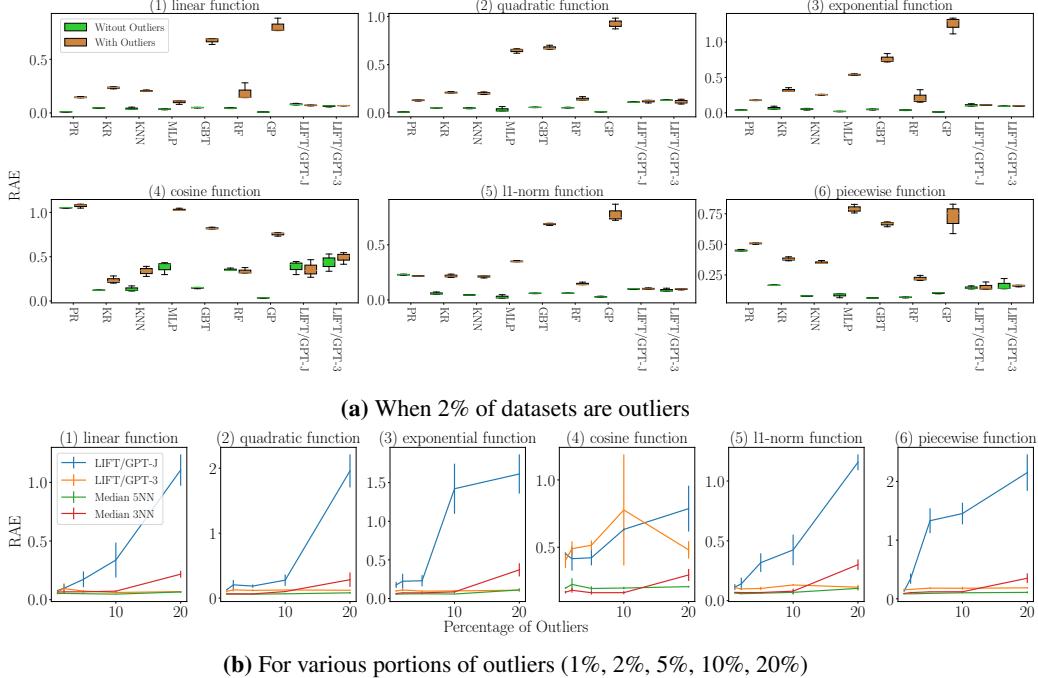


**(b)** Regression tasks (function approximation). The regression performance is reported in RMSE ( $\downarrow$ ).

**Figure 8: Sample complexity evaluations on classification and regression tasks.** Each figure presents the comparison of performance evaluated on LIFT/GPTs and baselines varying numbers of training samples (10–500 for classification and 10–1000 for regression). LIFT needs a slightly larger sample complexity to start achieving similar performances to the best baseline methods. For regression tasks, we note that LIFT achieves competitive or even better performance when around 1000s of samples are given, especially for the discontinuous functions, e.g., piecewise function.

Dataset (ID)	#Prompts	ODC	GPT-J			GPT-3		
			In-Context	LIFT/Subset	LIFT/Full-data	In-Context	LIFT/Subset	LIFT/Full-data
Breast (13)	35	70.69	56.90 $\pm$ 19.51	<b>58.62<math>\pm</math>2.44</b>	64.94 $\pm$ 11.97	62.07 $\pm$ 1.41	<b>70.69<math>\pm</math>0.00</b>	71.26 $\pm$ 1.62
TAE (48)	50	35.48	<b>34.33<math>\pm</math>1.47</b>	32.26 $\pm$ 9.50	61.29 $\pm$ 4.56	<b>37.64<math>\pm</math>4.02</b>	33.33 $\pm$ 1.52	65.59 $\pm$ 6.63
Vehicle (54)	14	25.88	<b>25.49<math>\pm</math>0.55</b>	26.04 $\pm$ 1.69	64.31 $\pm$ 2.37	<b>28.82<math>\pm</math>2.10</b>	23.73 $\pm$ 2.27	70.20 $\pm$ 2.73
Hamster (893)	43	53.33	48.89 $\pm$ 3.14	<b>60.00<math>\pm</math>10.88</b>	55.55 $\pm$ 16.63	<b>57.78<math>\pm</math>6.29</b>	53.33 $\pm$ 0.00	53.33 $\pm$ 0.00
Customers (1511)	29	68.18	56.06 $\pm$ 17.14	<b>59.85<math>\pm</math>2.84</b>	85.23 $\pm$ 1.61	60.61 $\pm$ 1.42	<b>63.26<math>\pm</math>6.96</b>	84.85 $\pm$ 1.42
LED (40496)	33	68.67	10.00 $\pm$ 0.82	<b>13.04<math>\pm</math>3.27</b>	65.33 $\pm$ 0.47	8.00 $\pm$ 1.63	<b>11.33<math>\pm</math>2.62</b>	69.33 $\pm$ 2.05

**Table 9: Comparison of accuracies ( $\uparrow$ ) between in-context learning and finetuning with LIFT on OpenML datasets.** “LIFT/Full-Data” and “LIFT/Subset” columns present results of LIFT on the full dataset and on the subset of data used in the corresponding in-context learning setting. We especially focus on the comparison between in-context learning and LIFT/subset. We put bold to the higher accuracy between the two methods. Note that the number of prompts for each dataset is selected to fit in the maximum context length that OpenAI GPT-3 can accept. We can see that LIFT/GPTs on full data achieve the best performances overall. When the same number of samples is used, applying LIFT outperforms the in-context learning in most cases for GPT-J, while the performances of LIFT are more comparable to those of in-context learning for GPT-3. Note that it is possible that in-context learning or LIFT/Subset performs worse than optimal deterministic classifier (ODC) since the number of training samples is limited.



**Figure 10: Comparing robustness of methods against outliers on regression tasks when the datasets contain (a) 2% outliers and (b) various portions of outliers.** We report each algorithm’s regression error measured by Related Absolute Error (RAE). (a) When training datasets contain 2% outliers, all LIFT models are highly robust against outliers compared with baselines. (b) When we increase the fraction of outliers (up to 20%), LIFT/GPT-3 is comparable to the strong baseline (median KNN), while LIFT/GPT-J fails.

outliers when there are 2% outliers in the training set. We observe that both LIFT/GPT-3 and LIFT/GPT-J are the most robust among all eight methods. To be specific, the performance of LIFT is almost unaffected, while all the baselines suffer from a huge performance drop. To further study how much LIFT is robust to outliers, in Fig. 10b, we evaluate the performance of LIFT on datasets including 1%, 2%, 5%, 10%, and 20% outliers, respectively. We include median-3NN and median-5NN, which are known to be robust to outliers [56] as baselines. Fig. 10b shows that LIFT with GPT-3 is almost as robust as median-3NN and median-5NN, while LIFT with GPT-J is more vulnerable when the percentage of outliers increases.

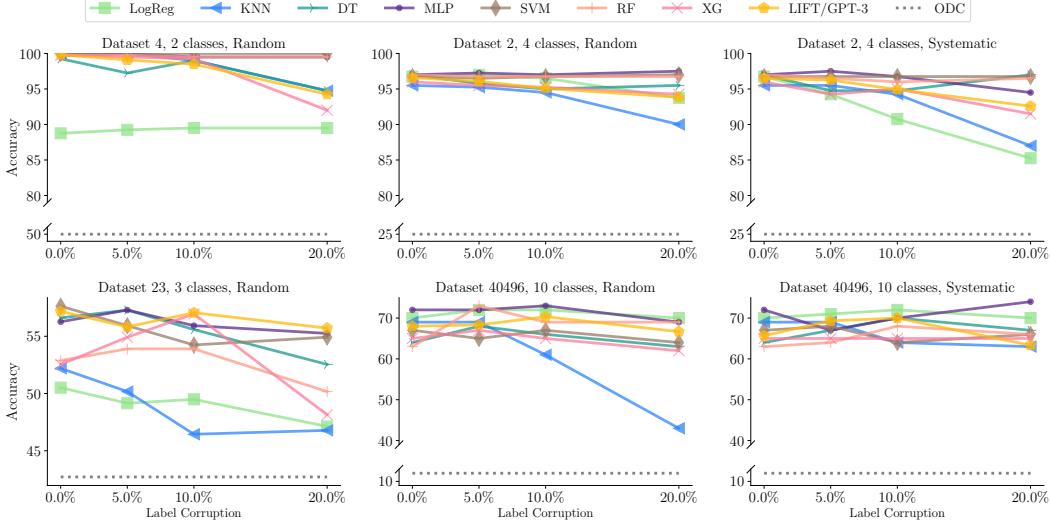
### 3.4.2 Robustness to Label Corruption on Training Data

We randomly corrupt some percentages of labeled samples in the training data by assigning different labels to the chosen samples. We apply two corruption schemes: random errors and systematic errors. For *random errors*, we randomly select a new label in the label space with an equal probability for all labels. For *systematic errors*, as described in [57], we replace a label with a determined label in the label space. In this experiment, we simply replace a label with its next label in the target label list, *e.g.*, 0 → 1, 1 → 2, 2 → 0 for a 3-class classification problem with three labels 0, 1, 2.

Fig. 11 presents our evaluations. We find that LIFT/GPT-3 can perform well in the presence of label corruption. In particular, LIFT follows the general trend of the evaluated algorithms, not exceptionally outperforming or underperforming in terms of robustness. We note that LIFT/GPT-3 almost always displays greater robustness to KNN.

### 3.4.3 Robustness to Class-Imbalance of Training Data

We evaluate LIFT on class-imbalanced classification tasks (OpenML datasets Pizza, Climate, and Customers having IDs 1444, 1467, and 1511), shown in Table 12. We use additional metrics: F1, precision, and recall (higher scores indicate better performance), which are considered better measurements for the imbalanced data than the accuracy. The higher values of the ODC’s accuracy



**Figure 11: Robustness against label corruption.** Each figure presents classification accuracies ( $\uparrow$ ) evaluated under different percentages of corruption in the training data (0% – 20%). We use synthetic data Blobs and Moons (ID 2 & 4) and real OpenML datasets CMC and LED (ID 23 & 40496). We simulate *random errors* (the first two columns) and *systematic errors* (the last column). LIFT/GPT-3 displays robustness across the datasets.

Dataset (ID)	Imb. Ratio		ODC	DC-0	DC-1	LogReg	KNN	DT	MLP	RBF-SVM	RF	XG	LIFT/GPT-3	LIFT/GPT-3
Pizza (1444)	7.36	Accuracy	<b>88.04</b>	<b>88.04</b>	11.96	86.92 $\pm$ 0.23	87.56 $\pm$ 0.68	87.24 $\pm$ 0.60	86.28 $\pm$ 1.37	<b>88.04<math>\pm</math>0.00</b>	<b>88.04<math>\pm</math>1.04</b>	<b>88.04<math>\pm</math>0.68</b>	83.89 $\pm$ 0.45	85.17 $\pm$ 1.35
		F1	0.00	0.00	21.37	10.77 $\pm$ 2.74	16.84 $\pm$ 5.26	9.01 $\pm$ 12.74	11.77 $\pm$ 4.43	0.00 $\pm$ 0.00	15.79 $\pm$ 6.45	35.50 $\pm$ 3.68	<b>35.83<math>\pm</math>3.61</b>	24.52 $\pm$ 1.78
		Precision	0.00	0.00	11.96	28.97 $\pm$ 3.41	42.86 $\pm$ 10.10	13.89 $\pm$ 19.64	32.69 $\pm$ 12.26	0.00 $\pm$ 0.00	51.67 $\pm$ 23.21	<b>52.21<math>\pm</math>7.28</b>	38.84 $\pm$ 5.24	32.41 $\pm$ 6.55
		Recall	0.00	0.00	<b>100.00</b>	6.67 $\pm$ 1.89	10.67 $\pm$ 3.77	6.67 $\pm$ 4.43	8.00 $\pm$ 3.27	0.00 $\pm$ 0.00	9.33 $\pm$ 3.77	28.00 $\pm$ 5.66	33.33 $\pm$ 2.36	20.00 $\pm$ 0.00
Climate (1467)	11.00	Accuracy	<b>91.67</b>	8.33	<b>91.67</b>	88.89 $\pm$ 0.76	90.74 $\pm$ 0.76	88.89 $\pm$ 2.27	<b>91.67<math>\pm</math>0.00</b>	87.96 $\pm$ 0.00	91.36 $\pm$ 0.44	89.51 $\pm$ 0.87	87.04 $\pm$ 2.27	<b>91.67<math>\pm</math>0.00</b>
		F1	<b>95.65</b>	0.00	<b>95.65</b>	94.00 $\pm$ 0.43	95.13 $\pm$ 0.49	94.04 $\pm$ 1.27	<b>95.65<math>\pm</math>0.00</b>	93.47 $\pm$ 0.00	95.48 $\pm$ 0.24	94.37 $\pm$ 0.48	94.51 $\pm$ 1.08	<b>95.65<math>\pm</math>1.00</b>
		Precision	91.67	0.00	91.67	<b>93.07<math>\pm</math>0.06</b>	91.85 $\pm$ 0.43	92.26 $\pm$ 1.20	91.67 $\pm$ 0.00	93.00 $\pm$ 0.00	91.64 $\pm$ 0.04	92.83 $\pm$ 0.43	92.10 $\pm$ 0.36	91.67 $\pm$ 0.00
		Recall	<b>100.00</b>	0.00	<b>100.00</b>	94.95 $\pm$ 0.82	98.65 $\pm$ 0.48	95.96 $\pm$ 2.86	<b>100.00<math>\pm</math>0.00</b>	93.94 $\pm$ 0.00	99.66 $\pm$ 0.48	95.96 $\pm$ 0.82	97.08 $\pm$ 2.57	<b>100.00<math>\pm</math>0.00</b>
Customers (1511)	2.14	Accuracy	68.18	68.18	31.82	87.12 $\pm$ 0.54	<b>88.64<math>\pm</math>0.00</b>	85.98 $\pm$ 0.53	86.36 $\pm$ 1.86	86.36 $\pm$ 0.00	85.23 $\pm$ 0.00	85.23 $\pm$ 0.00	85.23 $\pm$ 1.61	84.85 $\pm$ 1.42
		F1	0.00	0.00	48.28	79.76 $\pm$ 0.89	80.51 $\pm$ 0.36	78.60 $\pm$ 1.00	77.80 $\pm$ 2.79	78.82 $\pm$ 0.35	76.64 $\pm$ 0.39	76.91 $\pm$ 0.39	<b>84.43<math>\pm</math>1.43</b>	75.28 $\pm$ 2.60
		Precision	0.00	0.00	31.82	79.79 $\pm$ 1.23	<b>88.64<math>\pm</math>1.61</b>	76.40 $\pm$ 0.38	81.00 $\pm$ 4.65	77.94 $\pm$ 0.90	77.14 $\pm$ 0.91	76.50 $\pm$ 0.91	87.57 $\pm$ 7.11	78.18 $\pm$ 1.97
		Recall	0.00	0.00	100.00	79.76 $\pm$ 1.68	73.81 $\pm$ 1.68	80.95 $\pm$ 1.68	75.00 $\pm$ 2.91	79.76 $\pm$ 1.68	76.19 $\pm$ 1.68	77.38 $\pm$ 1.68	<b>82.61<math>\pm</math>7.10</b>	72.62 $\pm$ 3.37

**Table 12: Comparing accuracy ( $\uparrow$ ), F1 ( $\uparrow$ ), Precision ( $\uparrow$ ), and Recall ( $\uparrow$ ) on imbalanced datasets in OpenML (Pizza, Climate, Customers).** All datasets are for binary classification and are highly imbalanced. The class-imbalance ratio (Imb. Ratio) is defined as the ratio of the number of samples in the majority class and that in the minority class. Here, DC-0 and DC-1 refer to deterministic classifiers that constantly predicts all samples as class 0 and 1 respectively. ODC refers to the optimal deterministic classifier that returns the major class learned from the training dataset. LIFT/GPTs achieve comparably high scores across the three tasks. For instance, LIFT/GPT-J achieves the best F1 on datasets Pizza and Customers.

imply the higher levels of imbalance in the data (50% shows the perfect balance). For the reference, we report the performances of the deterministic classifiers that always return the label of class 0 (DC-0) and class 1 (DC-1).

Though evaluated datasets all have high class-imbalance ratios, we find that LIFT can perform well, achieving high F1, precision, and recall scores across the tasks. For instance, on the Customers dataset (the class-imbalance ratio is nearly 8), ODC gets 0 for both precision and recall as all predicted labels of ODC are 0 (*that is*, the major class), while LIFT/GPT-J achieves the best recall ( $82.61 \pm 7.10$ ) and F1 scores ( $84.43 \pm 1.43$ ). Here, the 0 value of precision and recall in ODC means that ODC classifies all samples as negative, which is the major class in the training dataset.

#### 3.4.4 Robustness to Feature Corruption on Test Data

We check whether LIFT trained on clean samples is robust against the test data with feature corruption. Given clean test data  $(\mathbf{x}, y)$  having feature  $\mathbf{x} \in \mathbb{R}^p$  and label  $y$ , we explore whether adding small perturbation  $\delta$  on the feature changes the performance of trained LIFT, by checking the accuracy of LIFT on perturbed data  $(\mathbf{x} + \delta, y)$ . Given a perturbation budget  $\varepsilon \geq 0$ , we consider two types of perturbation  $\delta$  with  $\|\delta\|_\infty \leq \varepsilon$ : random noise and adversarial perturbation [55]. For random noise  $\delta$ , we test on two types: (1) random Gaussian noise  $\delta \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_p)$  scaled to satisfy  $\|\delta\|_\infty = \varepsilon$ , and (2)

Source Target	Random noise (Gaussian)			Random noise (signed const.)			PGD attack on LeNet-5			PGD attack on MLP		
	LeNet-5	MLP	LIFT/GPT-3	LeNet-5	MLP	LIFT/GPT-3	LeNet-5	MLP	LIFT/GPT-3	LeNet-5	MLP	LIFT/GPT-3
$\varepsilon = 0$	99.22	98.09	98.15	99.22	98.09	98.15	99.22	98.09	98.15	99.22	98.09	98.15
$\varepsilon = 0.01$	99.25	98.05	98.28	99.26	98.08	88.05	97.27	97.77	44.88	99.15	96.89	44.46
$\varepsilon = 0.1$	99.20	97.70	88.38	99.06	97.39	68.80	26.80	93.99	33.66	96.98	23.12	23.62
$\varepsilon = 0.3$	98.01	87.69	54.80	79.80	74.20	29.68	0.00	36.62	20.31	41.51	0.00	20.29

**Table 13: Accuracies ( $\uparrow$ ) of LIFT and baselines (LeNet-5, MLP) under the perturbation on the input feature of MNIST data.** Given the perturbation budget  $\varepsilon \in [0, 1]$ , we test on four types of perturbations within  $L_\infty$  ball of radius  $\varepsilon$ . (1): adding random Gaussian noise that is scaled to reach the  $L_\infty$  ball, (2): adding *signed constant* noise vector where each element has magnitude  $\varepsilon$  and random sign, (3) & (4): adversarial examples generated from a source network (LeNet-5 & MLP, respectively) using PGD attack [60] from foolbox [61]. For small perturbation radii ( $\varepsilon = 0.01$ ), LIFT/GPT-3 maintains high accuracy for random noise, both for Gaussian and signed constant noise types. When  $\varepsilon = 0.01$  or  $\varepsilon = 0.1$ , the performance of LIFT/GPT-3 for random noise and transferred adversarial attacks have significant gap, showing that the adversarial examples generated at LeNet-5 and MLP are transferred to LIFT/GPT-3.

signed constant noise  $\delta$  where each element  $\delta_i$  has magnitude  $\varepsilon$  and random sign. For adversarial perturbation  $\delta$ , we tested on the transfer attack [58] since we do not have full access to the GPT-3 model and finding adversarial examples in the discrete input space is complex [59]. In other words, we generate an adversarial example  $(\mathbf{x} + \delta, y)$  for a source neural network (that we can access) with constraint  $\|\delta\| \leq \varepsilon$ , and test whether the target network correctly classifies the adversarial examples.

Table 13 shows the results of LIFT and baselines for the MNIST classification problem. We tested on random noise (Gaussian and signed constant) and PGD attacks transferred from LeNet-5 and MLP. Given MNIST images having pixel values within  $[0, 1]$ , the perturbation radius is set to  $\varepsilon \in [0, 0.01, 0.1, 0.3]$ . We compared the results for three networks: LeNet-5, MLP (having 2 hidden layers, each with 300 neurons and 100 neurons), and LIFT/GPT-3. It is shown that LIFT/GPT-3 tolerates random noise (both Gaussian and signed constant) for small perturbation radius  $\varepsilon = 0.01$ . For  $\varepsilon \in \{0.01, 0.1\}$  value, one can observe a huge gap between the accuracy of LIFT/GPT-3 on random noise and that on the transferred adversarial attack, implying that the adversarial attack on LeNet-5 and MLP are transferred to LIFT/GPT-3.

We do not include the result for LIFT/GPT-J since it is not even robust against simple noise. Please refer to Section 5.2 to check the vulnerability of LIFT/GPT-J against test-time noise and how data augmentation improves the robustness of LIFT/GPT-J.

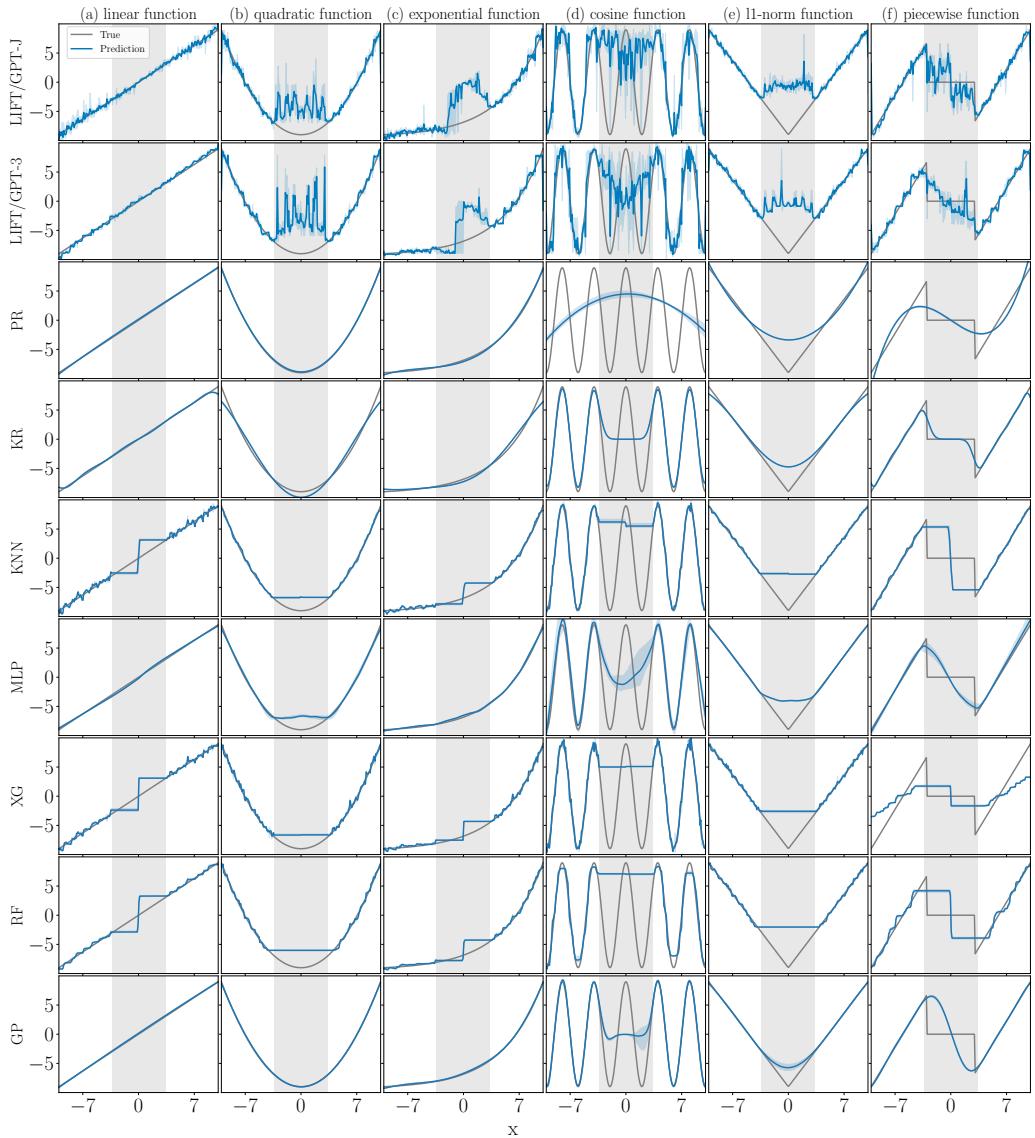
### 3.5 Can LIFT Interpolate and Extrapolate?

We investigate the interpolation and extrapolation performance of LIFT for regression tasks. Fig. 14 and Fig. 15 visualize the interpolation and extrapolation of various methods. All methods fail to extrapolate and interpolate well for all functions. It turns out that LIFT is not having good interpolation performance except in the linear regression case. An interesting observation is that LIFT tends to output seen values (from training data) for extrapolation. For example, in Fig. 15b, the outputs of LIFTs for  $x \notin [-10, 10]$  (extrapolation) lie in the range of outputs for  $x \in [-10, 10]$  (trained data), and similar behaviors are observed for other functions as well.

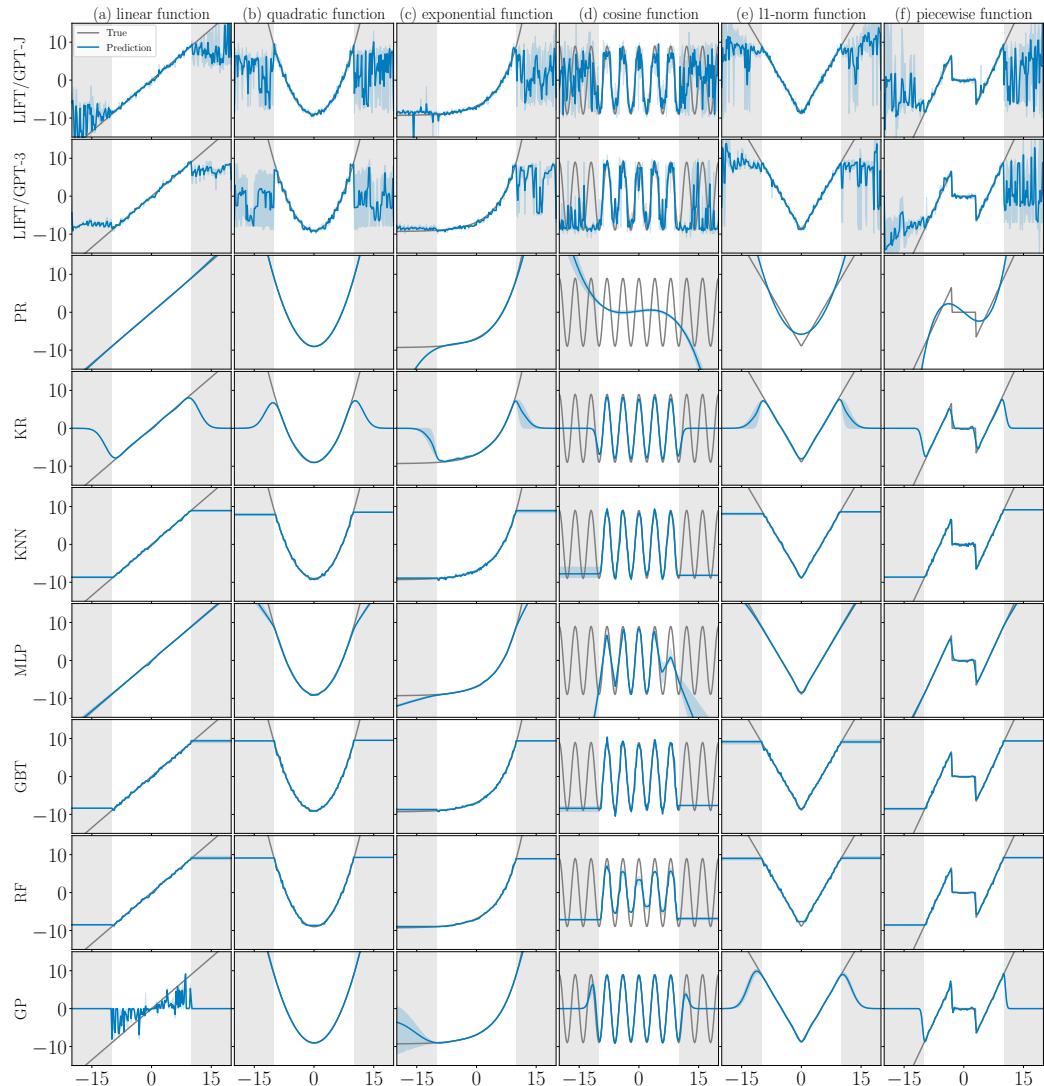
The appendix provides similar results for the regression problem with 2-dimensional input. Fig. 36 shows the results of the interpolation experiments when the number of training samples is  $n = 200$ .

### 3.6 Effects of Language Model Choices

In this experiment, we apply LIFT to different pretrained LMs to verify whether LIFT benefits more from larger LMs. Together with previously used GPT-J and GPT-3 (the version named Ada), we consider three bigger versions of GPT-3, namely Baggage, Curie, and Davinci (in the ascending order of the number of parameters). We compare all models on several classification tasks in Table 16 and regression tasks in Table 17. Overall, we find that the performance gain of using larger LMs is not consistently significant for LIFT. Although larger LMs outperform smaller LMs in many cases, the improvements are relatively small.



**Figure 14: Interpolation performance on synthetic regression tasks.** Each algorithm is trained with samples in the white background region ( $3 \leq |x| \leq 10$ ), and tested on the interpolation area  $|x| \leq 3$ . LIFT/GPTs are having worse interpolation performances compared with existing methods.



**Figure 15: Comparison of the extrapolation performance of LIFT and various baselines on synthetic regression tasks of approximating six functions  $f$ .** Each algorithm is trained by 200 samples  $(x, y)$  where the input  $x$  is drawn from interval  $[-10, 10]$  and the output is defined as  $y = f(x)$ . We test the how each algorithm perform regression for  $x \notin [-10, 10]$ .

**Verifying the capability of large LMs, when LIFT is not used.** We first verify if larger LMs are more helpful for the evaluated downstream tasks. We evaluate LMs in the in-context classification when no fine-tuning (LIFT) is involved. Table 16 shows consistent improvements in classification performance when the size of LMs increases across all the tasks. Thus, larger LMs, with larger embedded knowledge are more useful for these downstream tasks.

**When LIFT is used.** Both Table 16 and Table 17 show that using larger LMs may positively affect LIFT in several tasks and settings compared to the smaller LMs. However, the performance gains from replacing the smaller LMs with larger LMs are not consistent across the settings. For instance, in the classification settings without feature names, Davinci performs better than GPT-J on four datasets and worse than on one dataset. For the setting with feature names, GPT-J performs better than Davinci on two out of three tasks. Furthermore, the performance gains of large LMs over the smaller models are not relatively significant. We note that LIFT always outperforms the in-context learning using the same pretrained LMs in most cases. The regression results shown in Table 17 further confirm that the improvement from utilizing larger LMs is relatively small.

Tasks		LIFT Classification W/O Feat. Names					LIFT Classification W/ Feat. Names		
Dataset (ID)		Customers (1511)	Texture (1493)	Margin (1491)	TAE (48)	Vehicle (54)	TAE (48)	CMC (23)	Vehicle (54)
LIFT/GPT-J		93.97±1.00	50.32±2.18	50.23±1.33	61.29±6.97	64.31±2.37	67.74±11.48	48.36±0.97	69.02±3.67
LIFT/GPT-3	Ada	95.39±0.67	67.50±1.42	59.37±0.92	65.59±6.63	70.20±2.73	67.74±2.63	57.48±1.14	72.16±2.00
	Babbage	96.81±0.07	62.19±1.80	67.50±3.87	61.29±6.97	72.06±3.82	64.52±6.97	57.06±2.15	70.00±1.44
	Curie	95.21±0.06	62.50±0.97	61.88±1.48	66.67±6.09	74.27±0.73	65.59±4.02	55.42±0.84	70.66±2.28
	Davinci	96.81±0.41	57.19±0.70	58.13±2.50	64.52±9.50	71.47±0.88	65.59±6.63	56.31±0.04	68.16±1.69
Tasks		In-context Classification							
Dataset (ID) / #Prompts		TAE (48)/50	Breast (13)/35	LED (40496)/32	Customers (1511)/28	Vehicle (54)/42	Hamster (893)/13		
GPT-J		34.33±1.47	56.90±19.51	10.00±0.82	56.06±17.14	25.49±0.55	48.89±3.14		
GPT-3	Ada	37.64±4.02	62.07±1.41	8.00±1.63	60.61±1.42	28.82±2.10	57.78±6.29		
	Babbage	47.31±3.04	71.26±0.81	11.00±0.00	53.79±12.07	24.32±0.56	53.33±5.44		
	Curie	32.26±0.00	70.69±0.00	20.67±4.78	67.80±0.53	26.28±2.22	53.33±0.00		
	Davinci	49.46±4.02	67.82±4.06	20.67±6.60	68.94±0.54	26.28±2.22	55.55±3.14		

**Table 16: The effects of larger LMs under different classification settings.** Recall that our previous results on GPT-3 are based on the smallest model Ada. Here we use larger GPT-3 versions (Babbage, Curie, Davinci) as the pretrained LMs in our framework and evaluate the classification accuracy ( $\uparrow$ ) of them in three settings: classification *without* feature name, classification *with* feature name and in-context classification. For the setting of classification *with* feature names, we incorporate names of features (columns) into the input prompts (see more details in Sec. 4.1). For in-context learning, the OpenML dataset ID and number of prompts are written together at each column, *e.g.*, TAE (48)/50 means that we run experiments on the OpenML dataset TAE having ID 48, by using 50 input prompts. For the first two settings when LIFT is applied, larger LIFT/GPT-3 models (Babbage, Curie, Davinci) perform better than the smaller models LIFT/GPT-3-Ada and LIFT/GPT-J, but the performance gains are not always consistent and significant with model sizes. For the in-context classification (LIFT is not used), we observe more consistent improvement by using larger models.

Function \ Method	LIFT/GPT-J	LIFT/GPT-3			
		Ada	Babbage	Curie	Davinci
linear	0.08±0.01	0.06±0.01	0.06±0.00	0.06±0.01	0.06±0.00
quadratic	0.11±0.00	0.13±0.00	0.11±0.02	0.10±0.01	0.09±0.00
exponential	0.11±0.02	0.09±0.00	0.09±0.01	0.08±0.00	0.08±0.00
cosine	0.38±0.08	0.44±0.10	0.41±0.06	0.38±0.01	0.38±0.05
L1-norm	0.10±0.00	0.09±0.01	0.10±0.01	0.08±0.01	0.09±0.01
piecewise	0.15±0.01	0.17±0.05	0.15±0.02	0.15±0.01	0.14±0.01

**Table 17: Comparison of LIFT on different LMs across regression tasks.** The regression performance is measured by RAE ( $\downarrow$ ). In general, LIFT/GPT-3 with Davinci model performs the best, but the gaps to other models are not always significant.

## 4 Evaluation of LIFT-Specific Learning Properties

We study the learning characteristics and behaviors that are more specific to LIFT under the same setting as previous sections on non-language machine learning tasks. In particular, we investigate the followings: the effect of incorporating feature names to LIFT (Sec. 4.1) and the ability of LIFT in quantifying the predictive uncertainty (Sec. 4.2). Besides, we provide further analysis on other properties of LIFT in Appendix, including the effect of numerical types (Sec. ??) and whether LIFT can perform ridge regression via data augmentation (Sec. B.3).

### 4.1 Does LIFT Benefit from Incorporating Feature Names?

Unlike standard machine learning algorithms, LIFT can be provided with the context information by incorporating the feature names and task descriptions in the prompts. Intuitively, this incorporation may improve the sample complexity of LIFT as the prior knowledge already learned in the pretraining phase may help LIFT predict better. We empirically verify this intuition and show our results in Table 18 for several classification tasks using pretrained GPT-3 models. We provide further evaluations with GPT-J models (Table 32) and results on regression tasks (Table 33) in Appendix.

**Prompt templates.** We design five prompts templates to assess how incorporating feature names affects the performance of LIFT. For instance, consider a data sample “ $x = (\text{English speaker}, 23, 3, \text{summer}, 19)$ ,  $y = 3$ ” from TAE dataset where the feature names are “native speaker, instructor, course, semester, class size”, and the target attribute is teaching performance. We can incorporate the contextual information by either simply replacing the “ $x_i$ ” in the prompts with the corresponding feature names or converting this sample into a

Dataset (ID)	ODC	XGBoost	LIFT				
			W/o Names	Shuffled-Names I	Shuffled-Names II	Correct-Names I	Correct-Names II
CMC (23)	42.71	52.43 $\pm$ 0.42	<b>57.74<math>\pm</math>0.89</b>	56.27 $\pm$ 2.06	57.06 $\pm$ 4.24	57.40 $\pm$ 1.09	56.27 $\pm$ 2.22
TAE (48)	35.48	66.67 $\pm$ 8.05	65.59 $\pm$ 6.63	60.22 $\pm$ 6.72	64.52 $\pm$ 8.53	<b>69.89<math>\pm</math>9.31</b>	<b>69.89<math>\pm</math>6.72</b>
Vehicle (54)	25.88	73.14 $\pm$ 0.28	70.20 $\pm$ 2.73	70.20 $\pm$ 5.34	69.22 $\pm$ 2.72	<b>75.29<math>\pm</math>2.04*</b>	

**Table 18: The effect of using feature names on LIFT.** We compare classification accuracy ( $\uparrow$ ) of LIFT/GPT-3 when feature names provided in the target dataset *are and are not* incorporated to the prompts. We provide four versions of LIFT when feature names are correctly incorporated (Correct-Names columns) and when feature names are randomly shuffled (Shuffled-Names columns). We evaluate models on three OpenML datasets, including CMC (23), TAE (48), and Vehicle (54). We also compare our models with two baselines: the optimal deterministic classifier (ODC) and XGBoost. As a result, all LIFT models achieve better performance than ODC. Among the evaluated models, LIFTs with correct feature names achieve the best accuracies on both TAE and Vehicle datasets while achieving the comparable accuracies to the best model on the CMC dataset.  
 \*Two designs of the prompt format result in the same template for the Vehicle dataset.

coherent sentence. Meanwhile, we also investigate how shuffled feature names affect the performance of LIFT by designing the prompts accordingly. For illustration purposes, we provide the example of the five prompt templates as below.

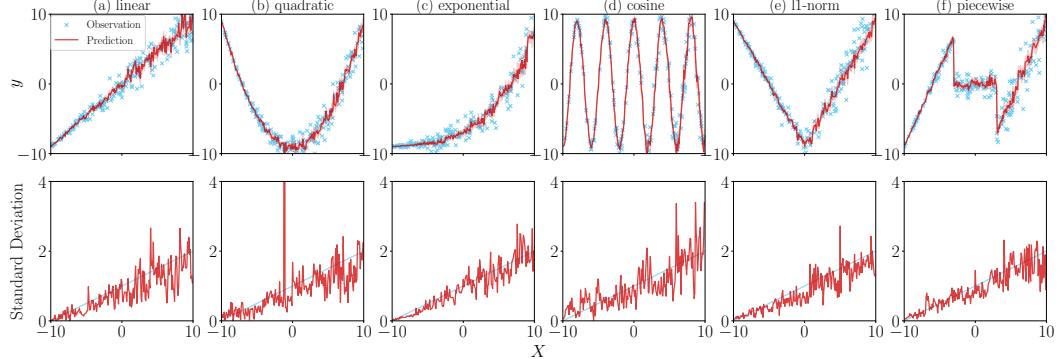
- (W/O Names) “When we have  $x_1 = 1, x_2 = 23, x_3 = 3, x_4 = 1, x_5 = 19$ , what should be  $y$  value?”
- (Correct-Names I) “When we have native speaker=English speaker, course instructor=23, course=3, semester=summer, class size=19, how is the teaching performance?”
- (Correct-Names II) “In the course 3 offered in the summer semester, there was a native English-speaking teaching assistant and an instructor whose ID is 23. How is the teaching performance?”
- (Shuffled-Names I) “When we have semester=English speaker, class size=23, semester=3, course instructor=summer, native speaker=19, how is the teaching performance?”
- (Shuffled-Names II) “In the course summer offered in the 3 semester, there was a 19 teaching assistant and an instructor whose ID is summer. How is the teaching performance?”

We note that the sentence generated using the (Shuffled-Names II) template can be incoherent.

**Settings.** (*Datasets*) Among the OpenML datasets evaluated in Table 4, we select three datasets: CMC, TAE, and Vehicle (with IDs being 23, 48, and 54) whose all provided feature names are meaningful and relevant to the prediction task and the response values. (*Baselines*) We compare our target model LIFT when feature names are correctly incorporated (Correct-Names I, II) with the versions of LIFT when feature names are incorrectly incorporated with randomly shuffled orders (Shuffled-Names I, II) and when feature names are not included (W/o Names). Also, we compare all models with the simple baseline ODC and the strong baseline XGBoost.

**Results.** Shown in Table 18 is our evaluation. At first, all LIFT models outperform ODC with large accuracy gaps, indicating that they are all properly trained. Overall, we observe that correctly incorporating feature names helps boost the performances of LIFT across the datasets. Moreover, on the OpenML datasets TAE and Vehicle, having IDs 48 and 54, using correct feature names helps improve the accuracy of the original LIFT to outperform XGBoost, which previously beat the versions of LIFT without feature names. For instance, using the correct feature names with the first template (Correct-Names I), LIFT improves accuracy on on Vehicle dataset from 70.20 to 75.29 to overcome XGBoost with an accuracy of 73.14. On the CMC dataset, the performance of LIFT remains almost similar after adding feature names, though all versions of LIFT outperform XGBoost with an accuracy gap of nearly 5%. Furthermore, if we use similar prompts but with shuffled feature names (Shuffled-Names I, II), then the performance of LIFT drops by a significant margin. Thus, this confirms that the aforementioned performance improvements are indeed due to proper prompting with correct feature/value association. These results validate that properly incorporating the feature names into LIFT may improve its prediction performance.

## 4.2 Can LIFT Quantify Predictive Uncertainty via Non-deterministic Decoding?



**Figure 19: Visualization of LIFT/GPT-J predictions under varying noise levels.** The training dataset (observations) for each function consists of 1000 samples  $\{(x_i, y_i)\}_{i=1}^{1000}$ , where  $y = f(x) + \sigma$ . The noise  $\sigma$  here monotonically increases along the  $x$ -axis. The predictions were performed on 103 evenly spaced samples in  $[-10, 10]$ . The standard deviations of the LIFT/GPT-J predictions were calculated based on 20 random predictions. As shown in the second row, the standard deviations of predictions from LIFT/GPT-J align well with that of noisy training samples (observations) across all datasets, implying that LIFT/GPT-J can properly quantify the predictive uncertainty.

Estimating the uncertainty of neural networks’ prediction has attracted a lot of interest [62, 63], with the majority of works involving Bayesian approaches [64, 62, 65]. Here, we investigate whether LIFT can properly quantify the predictive uncertainty by observing its behavior under various noise levels. Our results are shown in Fig. 19. The experiments are conducted on six synthetic regression datasets, each consisting of 1,000 noisy training samples shown as blue markers in the first row. To be specific, we generate (1) the input  $x$  following the guideline in Sec. A.1 for regression tasks and (2) the noisy outcome  $y$  where the standard deviation of noise  $\sigma(x) = (x + 10)/10$  increases along the  $x$ -axis (from  $x = -10$  to  $x = 10$ ), and study how different noise level affects the predictive behavior of LIFT. In the inference phase, we set the decoding temperature  $T = 1$  for LIFT to make random predictions. For visualization purposes, we generate an additional 103 samples uniformly in  $[-10, 10]$  for each task and plot the standard deviation of 20 LIFT/GPT-J predictions on each sample in the bottom row of Fig. 19. Note that the bottom row of Fig. 19 shows that the standard deviation of LIFT/GPT-J’s prediction nearly matches that of noisy training samples (observations) across different tasks. These results imply that LIFT/GPT-J can estimate the uncertainty of predictions. Similar behavior has been observed for LIFT/GPT-3, as shown in Fig. 31 of Sec. B.2.3.

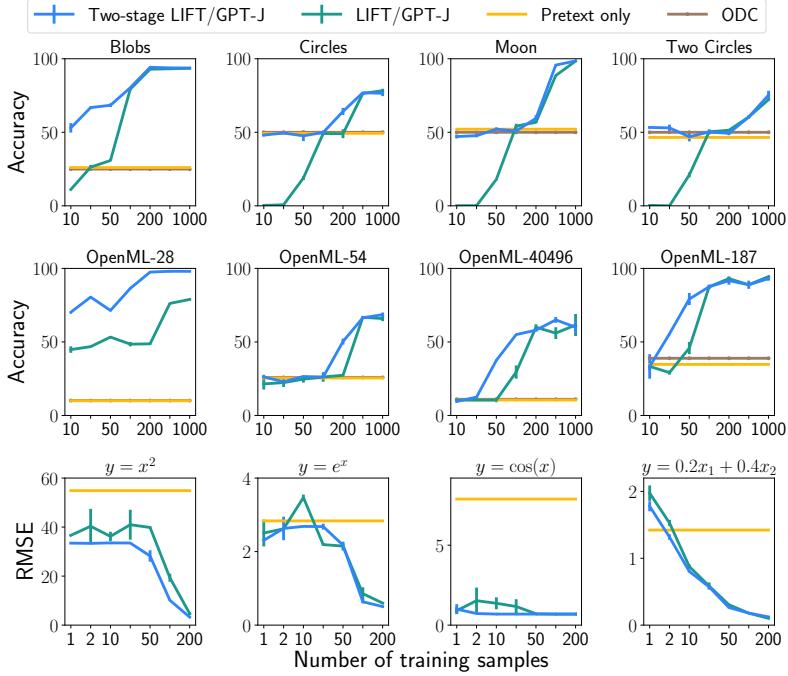
## 5 LIFT Combined with Advanced Techniques

In this section, we explore how well LIFT performs when combined with advanced techniques: two-stage fine-tuning (Sec. 5.1) and data augmentation (Sec. 5.2).

### 5.1 Two-stage Fine-Tuning

In Sec. 3.3, we observe that LMs need a sufficient number of samples to adapt to the non-language tasks. We suspect that this adaptation contains two phases. In the first phase, LMs use some samples to learn the task description, *i.e.*, input space, label space, and the sentence template [66, 54]. In the second phase, LMs use the rest of the samples to improve the performance on the target non-language task. Therefore, we consider utilizing synthetic data (instead of the actual training data) to describe the task for LMs in the first phase, thus reducing the sample complexity. This results in a new two-stage training procedure for LIFT.<sup>2</sup> Specifically, for any given dataset, we first generate two pretext tasks with simple synthetic Gaussian datasets (discussed in A.1) that share the same number of features and the label space (for classification tasks) or the range of responses’ values (for the

<sup>2</sup>We notice a recent relevant work [67] demonstrating the usefulness of the intermediate fine-tuning method for LMs. However, their focus is to propose self-supervised objectives for fine-tuning pretrained LMs in order to improve LMs’ few-shot in-context learning.



**Figure 20: Improving LIFT with the two-stage fine-tuning procedure.** Using GPT-J as the pretrained LM, we first apply LIFT on two pretext tasks with synthetic datasets before using the target dataset. We evaluate our methods on several classification tasks (first two rows, measured by accuracy ( $\uparrow$ )) and regression tasks (the last row, measured by RMSE). The OpenML datasets OPT (28), Vehicle (54), LED (40496), and Wine (187) are denoted by their OpenML ID. The two baselines are the optimal deterministic classifier (ODC, brown) and Pretext only (yellow). Here, Pretext only represents results of the LMs fine-tuned with LIFT only using pretext tasks. We can see that both versions of LIFT are better than two baselines. In most cases, when the number of training samples is small, the two-stage version of LIFT (blue) outperforms the original LIFT (green), especially on classification tasks.

regression tasks) to the actual data. We apply LIFT on pretext tasks for a few (2 or 3) epochs, then continue fine-tuning LIFT with the target (given) dataset. For GPT-3, it is unclear how to keep the order of samples not shuffled with the current black-box API during the fine-tuning stage. Hence, we only provide the experimental results of GPT-J.

Fig. 20 presents results from 12 datasets, including synthetic classification tasks (the top row), OpenML classification tasks (the second row), and regression tasks (the bottom row). We can see that the two-stage version of LIFT (blue line) achieves better performance compared to the vanilla version of LIFT (green line), especially when small numbers of training samples are provided. The effect of two-stage fine-tuning is more clearly shown in synthetic classification datasets (Blobs, Circles, Moon, and Two Circles). For instance, on the Blobs dataset, two-stage LIFT/GPT-J achieves more than 50% accuracy with only 10 samples, while the accuracy of LIFT/GPT-J is below the baselines (ODC and Pretext only). The two versions of LIFT are likely to converge when the number of samples increases. We emphasize that the purpose of pretexts is to inform and describe the task rather than reveal potential relationships between features. Therefore, generating the pretexts requires no knowledge about the correlation of the features, only the number of features and the number of classes.

## 5.2 Data Augmentation

Data augmentation [68] has been considered a simple tool for improving the generalization performance for various classification problems. In this section, we provide the preliminary result on the effect of data augmentation in LIFT. Table 21 shows the effect of adding random noise in the training data on the performance of LIFT/GPT-J for the MNIST classification problem. Here, we tested each model on three different settings: (1) clean data, (2) Gaussian noise, and (3) Signed constant

	Clean	Gaussian noise		Signed const. noise	
	$\varepsilon = 0$	$\varepsilon = 0.01$	$\varepsilon = 0.1$	$\varepsilon = 0.01$	$\varepsilon = 0.1$
LeNet-5	<b>99.22</b>	<b>99.25</b>	99.20	<b>99.26</b>	<b>99.06</b>
MLP	98.09	98.05	97.70	98.08	97.39
LIFT/GPT-J	<b>96.88</b>	<b>95.27</b>	56.14	55.83	27.73
LIFT/GPT-J, DA (Gaussian, $\varepsilon = 0.05$ )	93.80	94.39	93.40	93.46	61.24
LIFT/GPT-J, DA (Gaussian, $\varepsilon = 0.1$ )	93.78	94.31	<b>94.98</b>	<b>94.12</b>	<b>75.25</b>

**Table 21: Accuracies ( $\uparrow$ ) of LIFT with/without data augmentation (DA), as well as baselines (LeNet-5, MLP) on MNIST.** Each row represents different way of training, and each column means different test data. Here, data augmentation (DA) means that we are using a noisy version of MNIST training data by adding Gaussian noise. Given an MNIST image having range  $[0,1]$ , the noise is added in the  $L_\infty$  ball with radius  $\varepsilon$ . One can confirm that the data augmentation significantly improves the tolerance of LIFT/GPT-J against perturbed test data in both Gaussian and signed constant noise. For each column, we boldfaced the highest value among baselines and the highest value among LIFT/GPT-J.

noise. We allow each noise can perturb up to the magnitude of  $\varepsilon \in [0, 1]$  at each dimension (*i.e.*, each pixel) when the black/white pixel of MNIST is represented in the  $[0, 1]$  range. Following the design in Section 3.4.4, the noise  $\delta$  with  $\|\delta\|_\infty \leq \varepsilon$  is generated as below: the random Gaussian noise is defined as a scaled version of  $\mathcal{N}(0, \mathbf{I}_d)$  to satisfy with  $\|\delta\|_\infty = \varepsilon$ , whereas each element  $\delta_i$  of the signed constant noise has magnitude  $\varepsilon$  and random sign.

One can observe that LIFT/GPT-J without any data augmentation (DA) is vulnerable to random noise, unlike existing baselines (LeNet-5 and MLP). However, when we apply data augmentation, *i.e.*, train LIFT/GPT-J with noisy training data, the accuracy improves significantly for the perturbed (either adding Gaussian noise or Signed constant noise) test data. This shows the effectiveness of simple data augmentation in LIFT. Exploring the effect of other data augmentation schemes, *e.g.*, mixup [69] and its variants [70, 71, 72], is remained an interesting future work.

## 6 Related Works

**Pretraining and adapting language models (LMs).** Our work makes use of modern large LMs, which promoted significant advances in the field of natural language processing (NLP) [4]. Most popular LMs use the transformer architectures [12, 13] as the backbone, from early models like BERT [73] built on Transformer encoders to GPT variants [74, 21] built on Transformer decoders. Multiple modern large LMs have been proposed, including RoBERTa [75], ALBERT [76], XLNet [77], and the latest models with billions or trillions of parameters, such as GPT-3 [78], Switch-Transformers [79], and PALM [80].

LMs are trained to encode a large amount of linguistic knowledge from multiple sources in their contextual representations [81], which are useful and can be easily adapted to a variety of other tasks. Thus, starting with BERT [73], it has become a standard practice to pretrain and then fine-tune a large LM for plenty of downstream tasks in lieu of training a model from scratch for a specific task [73, 19, 4, 75, 82]. This technique greatly impacts a wide range of NLP tasks, such as language modeling [19, 80], question answering system [83, 84], text summarization [85], neural machine translation [86], and reasoning [87, 88].

However, the great performance of fine-tuned LMs without architecture changes has been mainly limited to NLP tasks so far. This work, instead, investigates whether we can leverage LM fine-tuning for non-language tasks across different modalities. Our work is highly motivated by Frozen Pretrained Transformer (FPT) [18], which directly adapts GPT-2 [21] pretrained on language tasks and textual data to other modalities. FPT freezes most pretrained parameters except the layer normalization layers and adds input and output layers for fine-tuning. The authors empirically show that GPT-2 can be efficiently fine-tuned for different modalities and domains, including vision and numeric computation. Nevertheless, FPT requires changes in the architecture and objective function to adapt to different data representations, while our method LIFT does not. Furthermore, we mainly focus on basic machine learning tasks, such as function approximation or tabular data classification.

Several works have attempted to extend the existing LMs to handle different types of input data, such as images [89, 90], audio [91], tabular data [92], and knowledge base [93] by updating the pretraining phase with these data and their corresponding tasks. For instance, XGPT [94] takes images as the input and uses the image captioning task as the primary task in the pretraining stage for GPT to generate images’ captions. Similarly, multiple works utilize pretrained LMs for generating text descriptions of images or videos in image captioning (VisualGPT [95]) and video captioning (VideoBERT [96], Unified VLP [97], UniVL [98]). SpeechBERT [91] also integrates LMs for speech recognition in the weakly-supervised setting to reduce the need for scarce supervised data. LMs can also adapt to numeric tasks [99] or other domains such as protein folding [100] or symbolic math solvers [101]. Recent works [102, 103, 92] pretrain LMs with large tabular datasets, improving the question answering systems by reasoning over the tables. Compared with these existing works, our work is unique in that it is based on GPT language models trained *only* on textual data.

**Analyzing the adaptability of LMs.** Similar to our work, recent works [20, 104] have made efforts to understand and quantify the feasibility and limitations of the adaptability of large LMs for upstream performance and downstream tasks. For instance, the recent work [20] built a benchmark of 500 small language tasks for testing the adaptability of LMs, observing that the LMs [105] can adapt well to an extensive range of complex tasks rather than just memorizing the learned patterns. BIG-bench [106] is recently introduced as a new benchmark for quantifying the capacity of LMs, consisting of more than 200 tasks on a diverse set of topics. Another relevant work [107] attempts to understand the effect of LM pretraining by studying how the transformer architecture, the backbone of LMs, succeeds at a designed synthetic task. Similar efforts in this line of work are to analyze the behaviors, representations, and inductive bias of pretrained LMs [108, 109, 52] or investigate different aspects of LMs [110, 111, 112]. For instance, a recent work [110], in the investigation of the difficulty of numeracy in LMs observes that transformer-based language models do not work well on complex numeric tasks and are sensitive to different formats of numeracy. Note that these existing works focus primarily on downstream language tasks, while we focus on adapting LMs on non-language tasks without any modification of the loss or architecture.

**Methods for adapting LMs.** The most common method for adapting LMs is *fine-tuning* [113] which aims to slightly adjust pretrained parameters for learning the downstream tasks [114, 115]. Fine-tuning can involve simple architecture modifications, such as adding linear layers [116, 117] or freezing parts of the network [18, 118, 119]. Fine-tuning can be improved with advanced techniques, such as multi-stage methods [120], intermediate fine-tuning [121, 122, 123], or self-supervised training [124]. The recent progress in fine-tuning LMs focuses on the *parameter-efficient* techniques for minimizing the number of fine-tunable parameters, including adapter-based fine-tuning [25, 26, 27] that adds and trains small residual blocks between transformer layers, freezing-based fine-tuning [125, 18, 126] that freezes most of the pretrained parameters and fine-tunes only tiny parts of the networks, and distillation-based fine-tuning [127]. In this line of work, LoRA [24] further reduces the number of trainable parameters in large LMs by approximating the weight updates using low-rank matrices without changing pretrained parameters. LoRA is used as the fine-tuning method for GPT-J in our LIFT/GPT-J framework. To directly adopt these fine-tuning methods used in LMs for non-language tasks, it is common practice to modify the input/output layers and the loss functions. However, these modifications might lead to undesired behaviors like catastrophic forgetting [113, 128]. On the other hand, our method LIFT uses the language interface for fine-tuning without making any changes to the architecture or the loss function.

In-context few-shot learning paradigm [129, 130, 28, 29, 131] suggests modifying only the inputs of LMs by adding a few examples of the downstream task. A critical part of these methods is reformulating the downstream task samples to the language modeling inputs [66, 54], resulting in multiple efforts in generating [132], searching [133], and properly tuning the prompts [134, 135, 136]. While these methods have shown great effectiveness for multiple NLP downstream tasks, it is unclear how to apply them to downstream tasks of other modalities. On the other hand, our work successfully adapts LMs for non-language tasks, further pushing the application boundaries of large-scale LMs.

**Deep learning for tabular datasets.** While deep neural networks have been successfully applied to various data types, such as images or text, they still face difficulties with a few classification and regression tasks on tabular data [137, 138], one of the most popular data types in practice. This may be due to the heterogeneous nature of tabular data, with their features being sparse, type-mixed, and

weaker in correlation than natural image-language data [139, 137]. Multiple deep learning methods and architectures have been proposed for tabular datasets, from making discrete decision trees more differentiable [140, 141], regularizing neural networks’ weights [142, 143], to recent attempts using attention-based modules [144, 145, 146]. Though these transformer-based models are the closest works to us in this line of work, their works focus on designing and improving architecture designs for specifically learning the tabular data rather than adapting the LMs. To the best of our knowledge, we are the first to thoroughly study large LM adaptation for tabular learning without architecture changes. Our work shows promising results of LMs in closing the gap to the best-performing methods, including tree-based ensemble algorithms (Random Forest [147] or XGBoost [148]).

**General-purpose models (generalist models).** A primary goal of our work is to push the limit of the existing generalist language models (*e.g.*, GPT-3 [19]) to other modalities and domains, supporting the idea of building a domain-and-modality agnostic generalist model. Early works [149, 14, 150] explored this idea by developing and training multi-task and multi-modal models on a wide range of diverse tasks to obtain better generalization and adaptation. The development of large-scale LMs has significantly contributed to the area of generalist models for languages [19, 151, 3], vision [152], visual language [153, 154, 155], and control problems [156]. These generalist models are usually trained with the scale on an extensive range of corpora, probably containing multiple modalities and domains. In this line of work, a general-purpose architecture [157] has also been studied to handle different input and output data types. Although LIFT primarily focuses on the LMs, it can be applied to other generalist models with LM-like architectures, such as GATO [156]. Furthermore, it is worth noticing that our work shares the general goal with automated machine learning (AutoML) [158, 159] in improving the usability of machine learning. AutoML automates the standard machine learning pipeline for model selection and hyperparameter tuning from a set of existing algorithms. At the same time, LIFT uses only a single pretrained LMs for solving all tasks.

## 7 Discussion

### 7.1 Limitations and Open Questions

Despite its promising performances on various tasks and settings, LIFT has some limitations, which open the interesting questions for further improvement.

**The difficulty of regression tasks.** For regression tasks, LIFT only performs well for approximating low-dimensional functions (see Fig. 1) and does not perform well for high-dimensional functions (see Table 28). Furthermore, some interesting phenomena observed in the classification tasks are not consistently observed in the regression tasks. For example, incorporating feature names in the prompts does not consistently improve LIFT in the regression tasks (see Sec. B.2.4).

As previously discussed in Sec. 3.1, the difficulty of regression tasks on LIFT may come from the classification loss function used in LMs. Due to the adoption of the classification loss function, two different predictions will lead to the same loss, even if one of the predictions is closer to the true  $y$  value. As a result, we also observe that a reduction in RAE does not necessarily imply a reduction in LM loss (see Fig. 26). Therefore, we use RAE as the criterion for model selection. Moreover, how LIFT understands numerical values may also limit the regression performance of LIFT. Recent works [160, 112, 161, 162] have illustrated the difficulty and failures of the existing LMs in understanding the numbers because two numbers with close values can have very different tokenizations [112]. Recent attempts [110, 163, 164, 165, 166] propose new encoding schemes of numbers to improve the LMs’ numerical capabilities, probably helping LIFT in the regression tasks.

A promising method for improving LIFT on regression is *level encoding*. The idea of level encoding is to discretize the continuous values of the output  $y$  to better utilize the classification loss of LMs. Assuming that the range of  $y$  is known, we can partition this range into a finite number of bins and represent all values in the same bin by a unique canonical representation in a way that the number of mismatched bits between the representations of two values is proportional to their absolute difference. For instance, for all real-value  $y \in [0, 3]$ , we can define three bins as  $\{[0, 1), [1, 2), [2, 3]\}$  with the canonical representations being 00, 01, 11. With these bins, 0.3 and 0.7 are represented as 00, and 1.5 and 1.1 are represented as 01. The distance between representations of 0.3 and 1.1 is only 1 bit, which is proportional to their absolute distance of 0.8. For the training of LIFT, we convert all

output values in the original training dataset into the level-encoding canonical representation and use them as the target values. By using the level encoding technique, the loss function of LMs can better capture the distance between the prediction and the true values, thus potentially improving the generalization of LIFT on regression tasks. We leave this as one of the interesting directions for our future investigation.

**The limitation of LIFT on classification tasks.** We observe that LIFT does not perform relatively well on classification tasks when the number of classes is large. For instance, Table 4 shows that the accuracies of LIFT/GPT-3 are lower than RBF-SVM and XGboost on the datasets with 100 classes. Another limitation is that the dimension of features LIFT can handle is upper bounded due to the limited context length of LMs. This limitation may be mitigated by using LMs with a more memory-efficient variant or implementation of transformer models, *e.g.*, see [167].

**Other open questions.** In addition to previously discussed questions of improving LIFT for regression and classification tasks, our pioneering work on LIFT is also expected to open up interesting research questions on generalist models. First, can generalist LMs (*e.g.*, GPTs) play a leading role in developing universal models that can adapt well to any data? Second, can we apply LIFT to different generalist models, such as GATO [156]?

## 7.2 Broader Impact

LIFT greatly simplifies the machine learning pipeline that requires only the reformatting of training datasets of the target task. This simplicity helps enable *no-code ML* for the masses, where general users without prior knowledge of the ML frameworks can use LIFT for their target non-language tasks by properly designing the input/output prompt format. Therefore, LIFT can apply to a wide range of applications and areas, such as credit loaning, disease diagnosis, and criminal sentencing. This is closely related to the line of automated machine learning research [158, 159], which aims to automate the standard machine learning methods pipeline.

Employing LIFT without careful justification or understanding will lead to undesired outcomes, such as discrimination. Since most existing language models (LMs) are pretrained on a large amount of human-annotated data, LIFT could exhibit discrimination against different demographic groups (*e.g.*, gender, race, ethnicity) due to the bias existing in the training datasets. In other words, LIFT may prefer certain groups while making decisions in downstream tasks, especially when feature names and different demographic contexts are fed at training and inference time. This effect is exacerbated by the use of large pretrained LMs (*i.e.*, GPT-J and GPT-3), which have been known to inherently contain bias [168]. The bias in the pretraining data for these large language models adds an opaque layer to regression and classification tasks beyond bias within the downstream data. Therefore, adopting LIFT in tasks that consider demographic information requires more consideration to avoid discrimination. To further remove the bias, users can combine LIFT with the existing fairness-aware reweighting mechanisms [169, 170] or data augmentation and parameter efficient fine-tuning techniques [171].

Finally, we emphasize that more model evaluation steps are required when applying LIFT instead of using it as a panacea for all applications. We believe our work can significantly benefit society by providing a simple tool for handling various tasks with proper justification.

## 8 Conclusion

We study the Language-Interfaced Fine-Tuning (LIFT) framework for solving non-language downstream tasks without altering the model architecture or loss function. LIFT first converts labeled samples into sentences and then fine-tunes pretrained LMs on the obtained sentences. Via an extensive empirical study on a wide range of non-language tasks, we find that LIFT performs relatively well on a variety of low-dimensional classification and regression tasks, nearly matching the performance of the best baseline models in most cases. Furthermore, LIFT can achieve better classification performances with the proper use of in-context feature names. We further improve LIFT with the two-stage fine-tuning and the data augmentation techniques. It is worth noting that our work is the first to thoroughly study the efficacy of fine-tuning pretrained language models on various non-language tasks *without* any architecture changes, which paves the road for enabling “no-code machine learning with language models.”

## References

- [1] David Forsyth and Jean Ponce. *Computer vision: A modern approach*. Prentice hall, 2011.
- [2] Alexander Kolesnikov, Alexey Dosovitskiy, Dirk Weissenborn, Georg Heigold, Jakob Uszkoreit, Lucas Beyer, Matthias Minderer, Mostafa Dehghani, Neil Houlsby, Sylvain Gelly, Thomas Unterthiner, and Xiaohua Zhai. An image is worth 16x16 words: Transformers for image recognition at scale. 2021.
- [3] KR1442 Chowdhary. Natural language processing. *Fundamentals of artificial intelligence*, pages 603–649, 2020.
- [4] Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, et al. On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*, 2021.
- [5] Fei-Yue Wang, Jun Jason Zhang, Xinhua Zheng, Xiao Wang, Yong Yuan, Xiaoxiao Dai, Jie Zhang, and Liuqing Yang. Where does alphago go: From church-turing thesis to alphago thesis and beyond. *IEEE/CAA Journal of Automatica Sinica*, 3(2):113–120, 2016.
- [6] Kai Arulkumaran, Antoine Cully, and Julian Togelius. Alphastar: An evolutionary computation perspective. In *Proceedings of the genetic and evolutionary computation conference companion*, pages 314–315, 2019.
- [7] Shih-Chung B Lo, Heang-Ping Chan, Jyh-Shyan Lin, Huai Li, Matthew T Freedman, and Seong K Mun. Artificial convolution neural network for medical image pattern recognition. *Neural networks*, 8(7-8):1201–1214, 1995.
- [8] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [9] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 2961–2969, 2017.
- [10] David E Rumelhart, Geoffrey E Hinton, and Ronald J Williams. Learning internal representations by error propagation. Technical report, California Univ San Diego La Jolla Inst for Cognitive Science, 1985.
- [11] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.
- [12] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- [13] David So, Quoc Le, and Chen Liang. The evolved transformer. In *International Conference on Machine Learning*, pages 5877–5886. PMLR, 2019.
- [14] Lukasz Kaiser, Aidan N Gomez, Noam Shazeer, Ashish Vaswani, Niki Parmar, Llion Jones, and Jakob Uszkoreit. One model to learn them all. *arXiv preprint arXiv:1706.05137*, 2017.
- [15] Karl Weiss, Taghi M Khoshgoftaar, and DingDing Wang. A survey of transfer learning. *Journal of Big data*, 3(1):1–40, 2016.
- [16] Dhanesh Ramachandram and Graham W Taylor. Deep multimodal learning: A survey on recent advances and trends. *IEEE signal processing magazine*, 34(6):96–108, 2017.
- [17] Ronan Collobert and Jason Weston. A unified architecture for natural language processing: Deep neural networks with multitask learning. In *Proceedings of the 25th international conference on Machine learning*, pages 160–167, 2008.
- [18] Kevin Lu, Aditya Grover, Pieter Abbeel, and Igor Mordatch. Pretrained transformers as universal computation engines. *arXiv preprint arXiv:2103.05247*, 2021.

- [19] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- [20] Belinda Z Li, Jane Yu, Madian Khabsa, Luke Zettlemoyer, Alon Halevy, and Jacob Andreas. Quantifying adaptability in pre-trained language models with 500 tasks. *arXiv preprint arXiv:2112.03204*, 2021.
- [21] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- [22] Dan Su, Yan Xu, Genta Indra Winata, Peng Xu, Hyeondey Kim, Zihan Liu, and Pascale Fung. Generalizing question answering system with pre-trained language model fine-tuning. In *Proceedings of the 2nd Workshop on Machine Reading for Question Answering*, pages 203–211, Hong Kong, China, November 2019. Association for Computational Linguistics.
- [23] Xuhui Zhou, Yue Zhang, Leyang Cui, and Dandan Huang. Evaluating commonsense in pre-trained language models. *CoRR*, abs/1911.11931, 2019.
- [24] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021.
- [25] Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. Parameter-efficient transfer learning for nlp. In *International Conference on Machine Learning*, pages 2790–2799. PMLR, 2019.
- [26] Sylvestre-Alvise Rebuffi, Hakan Bilen, and Andrea Vedaldi. Learning multiple visual domains with residual adapters. *Advances in neural information processing systems*, 30, 2017.
- [27] Ruize Wang, Duyu Tang, Nan Duan, Zhongyu Wei, Xuanjing Huang, Guihong Cao, Dixin Jiang, Ming Zhou, et al. K-adapter: Infusing knowledge into pre-trained models with adapters. *arXiv preprint arXiv:2002.01808*, 2020.
- [28] Yao Lu, Max Bartolo, Alastair Moore, Sebastian Riedel, and Pontus Stenetorp. Fantastically ordered prompts and where to find them: Overcoming few-shot prompt order sensitivity. *arXiv preprint arXiv:2104.08786*, 2021.
- [29] Bonan Min, Hayley Ross, Elior Sulem, Amir Pouran Ben Veyseh, Thien Huu Nguyen, Oscar Sainz, Eneko Agirre, Ilana Heinz, and Dan Roth. Recent advances in natural language processing via large pre-trained language models: A survey. *arXiv preprint arXiv:2111.01243*, 2021.
- [30] Kyle M. Monahan. Iris dataset for machine learning, 2020.
- [31] Ben Wang and Aran Komatsuzaki. GPT-J-6B: A 6 Billion Parameter Autoregressive Language Model. <https://github.com/kingoflolz/mesh-transformer-jax>, May 2021.
- [32] Openai fine-tuning documentation: Preparing your dataset. <https://beta.openai.com/docs/guides/fine-tuning/preparing-your-dataset>.
- [33] Daphne Ippolito, Reno Kriz, Maria Kustikova, João Sedoc, and Chris Callison-Burch. Comparison of diverse decoding methods from conditional language models. *arXiv preprint arXiv:1906.06362*, 2019.
- [34] Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. The curious case of neural text degeneration. *arXiv preprint arXiv:1904.09751*, 2019.
- [35] Openai fine-tuning documentation: Create completion. <https://beta.openai.com/docs/api-reference/completions/create>.
- [36] Joaquin Vanschoren, Jan N. van Rijn, Bernd Bischl, and Luis Torgo. Openml: Networked science in machine learning. *SIGKDD Explorations*, 15(2):49–60, 2013.

- [37] Yann LeCun. The mnist database of handwritten digits. <http://yann.lecun.com/exdb/mnist/>, 1998.
- [38] Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.
- [39] Ian J Goodfellow, Mehdi Mirza, Da Xiao, Aaron Courville, and Yoshua Bengio. An empirical investigation of catastrophic forgetting in gradient-based neural networks. *arXiv preprint arXiv:1312.6211*, 2013.
- [40] Keyulu Xu, Mozhi Zhang, Jingling Li, Simon Shaolei Du, Ken-Ichi Kawarabayashi, and Stefanie Jegelka. How neural networks extrapolate: From feedforward to graph neural networks. In *International Conference on Learning Representations*, 2021.
- [41] Medical insurance dataset (kaggle). <https://www.kaggle.com/datasets/mirichoi0218/insurance>.
- [42] Pinar Tüfekci. Prediction of full load electrical power output of a base load operated combined cycle power plant using machine learning methods. *International Journal of Electrical Power & Energy Systems*, 60:126–140, 2014.
- [43] John R Quinlan et al. Learning with continuous classes. In *5th Australian joint conference on artificial intelligence*, volume 92, pages 343–348. World Scientific, 1992.
- [44] Paulo Cortez and Alice Maria Gonçalves Silva. Using data mining to predict secondary school student performance. 2008.
- [45] Shai Shalev-Shwartz and Shai Ben-David. *Understanding machine learning: From theory to algorithms*. Cambridge university press, 2014.
- [46] Vadim Borisov, Tobias Leemann, Kathrin Seßler, Johannes Haug, Martin Pawelczyk, and Gjergji Kasneci. Deep neural networks and tabular data: A survey. October 2021.
- [47] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- [48] Benteng Ma, Xiang Li, Yong Xia, and Yanning Zhang. Autonomous deep learning: A genetic dcnn designer for image classification. *Neurocomputing*, 379:152–161, 2020.
- [49] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. *Advances in neural information processing systems*, 25, 2012.
- [50] Sara Sabour, Nicholas Frosst, and Geoffrey E Hinton. Dynamic routing between capsules. *Advances in neural information processing systems*, 30, 2017.
- [51] Cristian Ivan. Convolutional neural networks on randomized data. In *CVPR Workshops*, pages 1–8, 2019.
- [52] Charles Lovering, Rohan Jha, Tal Linzen, and Ellie Pavlick. Predicting inductive biases of pre-trained models. In *International Conference on Learning Representations*, 2020.
- [53] Gowthami Somepalli, Liam Fowl, Arpit Bansal, Ping Yeh-Chiang, Yehuda Dar, Richard Baraniuk, Micah Goldblum, and Tom Goldstein. Can neural nets learn the same model twice? investigating reproducibility and double descent from the decision boundary perspective. *arXiv preprint arXiv:2203.08124*, 2022.
- [54] Sewon Min, Xinxi Lyu, Ari Holtzman, Mikel Artetxe, Mike Lewis, Hannaneh Hajishirzi, and Luke Zettlemoyer. Rethinking the role of demonstrations: What makes in-context learning work? *arXiv preprint arXiv:2202.12837*, 2022.
- [55] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199*, 2013.

- [56] Peter J Huber. Robust statistics. In *International encyclopedia of statistical science*, pages 1248–1251. Springer, 2011.
- [57] Yanyao Shen and Sujay Sanghavi. Learning with bad training data via iterative trimmed loss minimization. In *International Conference on Machine Learning*, pages 5739–5748. PMLR, 2019.
- [58] Alexey Kurakin, Ian Goodfellow, and Samy Bengio. Adversarial machine learning at scale. *arXiv preprint arXiv:1611.01236*, 2016.
- [59] Wei Emma Zhang, Quan Z Sheng, Ahoud Alhazmi, and Chenliang Li. Adversarial attacks on deep-learning models in natural language processing: A survey. *ACM Transactions on Intelligent Systems and Technology (TIST)*, 11(3):1–41, 2020.
- [60] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. *arXiv preprint arXiv:1706.06083*, 2017.
- [61] Jonas Rauber, Wieland Brendel, and Matthias Bethge. Foolbox: A python toolbox to benchmark the robustness of machine learning models. *arXiv preprint arXiv:1707.04131*, 2017.
- [62] Yarin Gal and Zoubin Ghahramani. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In *international conference on machine learning*, pages 1050–1059. PMLR, 2016.
- [63] Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive uncertainty estimation using deep ensembles. *Advances in neural information processing systems*, 30, 2017.
- [64] Charles Blundell, Julien Cornebise, Koray Kavukcuoglu, and Daan Wierstra. Weight uncertainty in neural network. In *International conference on machine learning*, pages 1613–1622. PMLR, 2015.
- [65] Christos Louizos and Max Welling. Structured and efficient variational deep learning with matrix gaussian posteriors. In *International conference on machine learning*, pages 1708–1716. PMLR, 2016.
- [66] Laria Reynolds and Kyle McDonell. Prompt programming for large language models: Beyond the few-shot paradigm. In *Extended Abstracts of the 2021 CHI Conference on Human Factors in Computing Systems*, pages 1–7, 2021.
- [67] Mingda Chen, Jingfei Du, Ramakanth Pasunuru, Todor Mihaylov, Srinivas Iyer, Veselin Stoyanov, and Zornitsa Kozareva. Improving in-context few-shot learning via self-supervised training. *arXiv preprint arXiv:2205.01703*, 2022.
- [68] Connor Shorten and Taghi M Khoshgoftaar. A survey on image data augmentation for deep learning. *Journal of big data*, 6(1):1–48, 2019.
- [69] Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. mixup: Beyond empirical risk minimization. *arXiv preprint arXiv:1710.09412*, 2017.
- [70] Sangdoo Yun, Dongyoon Han, Seong Joon Oh, Sanghyuk Chun, Junsuk Choe, and Youngjoon Yoo. Cutmix: Regularization strategy to train strong classifiers with localizable features. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 6023–6032, 2019.
- [71] Jang-Hyun Kim, Wonho Choo, and Hyun Oh Song. Puzzle mix: Exploiting saliency and local statistics for optimal mixup. In *International Conference on Machine Learning*, pages 5275–5285. PMLR, 2020.
- [72] Jy-yong Sohn, Liang Shang, Hongxu Chen, Jaekyun Moon, Dimitris Papailiopoulos, and Kangwook Lee. Genlabel: Mixup relabeling using generative models. *arXiv preprint arXiv:2201.02354*, 2022.

- [73] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- [74] Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. Improving language understanding by generative pre-training. 2018.
- [75] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*, 2019.
- [76] Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. Albert: A lite bert for self-supervised learning of language representations. *arXiv preprint arXiv:1909.11942*, 2019.
- [77] Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Carbonell, Russ R Salakhutdinov, and Quoc V Le. Xlnet: Generalized autoregressive pretraining for language understanding. *Advances in neural information processing systems*, 32, 2019.
- [78] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- [79] William Fedus, Barret Zoph, and Noam Shazeer. Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity. *arXiv preprint arXiv:2101.03961*, 2021.
- [80] Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm: Scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311*, 2022.
- [81] Junyi Li, Tianyi Tang, Wayne Xin Zhao, Jian-Yun Nie, and Ji-Rong Wen. A survey of pretrained language models based text generation. *arXiv preprint arXiv:2201.05273*, 2022.
- [82] Shuang Li, Xavier Puig, Yilun Du, Clinton Wang, Ekin Akyurek, Antonio Torralba, Jacob Andreas, and Igor Mordatch. Pre-trained language models for interactive decision-making. *arXiv preprint arXiv:2202.01771*, 2022.
- [83] Bryan McCann, Nitish Shirish Keskar, Caiming Xiong, and Richard Socher. The natural language decathlon: Multitask learning as question answering. *arXiv preprint arXiv:1806.08730*, 2018.
- [84] Dan Su, Yan Xu, Genta Indra Winata, Peng Xu, Hyeondey Kim, Zihan Liu, and Pascale Fung. Generalizing question answering system with pre-trained language model fine-tuning. In *Proceedings of the 2nd Workshop on Machine Reading for Question Answering*, pages 203–211, 2019.
- [85] Dmitrii Aksenov, Julián Moreno-Schneider, Peter Bourgonje, Robert Schwarzenberg, Leonhard Hennig, and Georg Rehm. Abstractive text summarization based on language model conditioning and locality modeling. *arXiv preprint arXiv:2003.13027*, 2020.
- [86] Inigo Jauregi Unanue and Massimo Piccardi. Pretrained language models and backtranslation for English-Basque biomedical neural machine translation. In *Proceedings of the Fifth Conference on Machine Translation*, pages 826–832, Online, November 2020. Association for Computational Linguistics.
- [87] Xiao Liu, Yanan Zheng, Zhengxiao Du, Ming Ding, Yujie Qian, Zhilin Yang, and Jie Tang. Gpt understands, too. *arXiv preprint arXiv:2103.10385*, 2021.
- [88] Lya Hulliyatus Suadaa, Hidetaka Kamigaito, Kotaro Funakoshi, Manabu Okumura, and Hiroya Takamura. Towards table-to-text generation with numerical reasoning. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 1451–1465, 2021.

- [89] Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. *Advances in neural information processing systems*, 32, 2019.
- [90] Maria Tsimpoukelli, Jacob L Menick, Serkan Cabi, SM Eslami, Oriol Vinyals, and Felix Hill. Multimodal few-shot learning with frozen language models. *Advances in Neural Information Processing Systems*, 34:200–212, 2021.
- [91] Yung-Sung Chuang, Chi-Liang Liu, and Hung-Yi Lee. Speechbert: Cross-modal pre-trained language model for end-to-end spoken question answering. 2019.
- [92] Qian Liu, Bei Chen, Jiaqi Guo, Zeqi Lin, and Jian-guang Lou. Tapex: table pre-training via learning a neural sql executor. *arXiv preprint arXiv:2107.07653*, 2021.
- [93] Oshin Agarwal, Heming Ge, Siamak Shakeri, and Rami Al-Rfou. Knowledge graph based synthetic corpus generation for knowledge-enhanced language model pre-training. *arXiv preprint arXiv:2010.12688*, 2020.
- [94] Qiaolin Xia, Haoyang Huang, Nan Duan, Dongdong Zhang, Lei Ji, Zhifang Sui, Edward Cui, Taroon Bharti, and Ming Zhou. Xgpt: Cross-modal generative pre-training for image captioning. In *CCF International Conference on Natural Language Processing and Chinese Computing*, pages 786–797. Springer, 2021.
- [95] Jun Chen, Han Guo, Kai Yi, Boyang Li, and Mohamed Elhoseiny. Visualgpt: Data-efficient adaptation of pretrained language models for image captioning. *arXiv preprint arXiv:2102.10407*, 2021.
- [96] Chen Sun, Austin Myers, Carl Vondrick, Kevin Murphy, and Cordelia Schmid. Videobert: A joint model for video and language representation learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7464–7473, 2019.
- [97] Luwei Zhou, Hamid Palangi, Lei Zhang, Houdong Hu, Jason Corso, and Jianfeng Gao. Unified vision-language pre-training for image captioning and vqa. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 13041–13049, 2020.
- [98] Huaishao Luo, Lei Ji, Botian Shi, Haoyang Huang, Nan Duan, Tianrui Li, Jason Li, Taroon Bharti, and Ming Zhou. Univl: A unified video and language pre-training model for multimodal understanding and generation. *arXiv preprint arXiv:2002.06353*, 2020.
- [99] Jannis Born and Matteo Manica. Regression transformer: Concurrent conditional generation and regression by blending numerical and textual tokens. *arXiv preprint arXiv:2202.01338*, 2022.
- [100] John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf Ronneberger, Kathryn Tunyasuvunakool, Russ Bates, Augustin Žídek, Anna Potapenko, et al. Highly accurate protein structure prediction with alphafold. *Nature*, 596(7873):583–589, 2021.
- [101] Kimia Noorbakhsh, Modar Sulaiman, Mahdi Sharifi, Kallol Roy, and Pooyan Jamshidi. Pretrained language models are symbolic mathematics solvers too! *arXiv preprint arXiv:2110.03501*, 2021.
- [102] Jonathan Herzig, Paweł Krzysztof Nowak, Thomas Müller, Francesco Piccinno, and Julian Eisenschlos. TaPas: Weakly supervised table parsing via pre-training. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4320–4333, Online, July 2020. Association for Computational Linguistics.
- [103] Hiroshi Iida, Dung Thai, Varun Manjunatha, and Mohit Iyyer. Tabbie: Pretrained representations of tabular data. *arXiv preprint arXiv:2105.02584*, 2021.
- [104] Samira Abnar, Mostafa Dehghani, Behnam Neyshabur, and Hanie Sedghi. Exploring the limits of large scale pre-training. *arXiv preprint arXiv:2110.02095*, 2021.
- [105] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 21(140):1–67, 2020.

- [106] Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch, Adam R. Brown, Adam Santoro, et al. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models, 2022.
- [107] Yi Zhang, Arturs Backurs, Sébastien Bubeck, Ronen Eldan, Suriya Gunasekar, and Tal Wagner. Unveiling transformers with lego: a synthetic reasoning task, 2022.
- [108] Dieuwke Hupkes, Sara Veldhoen, and Willem Zuidema. Visualisation and diagnostic classifiers' reveal how recurrent and recursive neural networks process hierarchical structure. *Journal of Artificial Intelligence Research*, 61:907–926, 2018.
- [109] Alex Warstadt, Yian Zhang, Haau-Sing Li, Haokun Liu, and Samuel R Bowman. Learning which features matter: Roberta acquires a preference for linguistic generalizations (eventually). *arXiv preprint arXiv:2010.05358*, 2020.
- [110] Rodrigo Nogueira, Zhiying Jiang, and Jimmy Lin. Investigating the limitations of transformers with simple arithmetic tasks. *arXiv preprint arXiv:2102.13019*, 2021.
- [111] Guillaume Lample and François Charton. Deep learning for symbolic mathematics. In *International Conference on Learning Representations*, 2019.
- [112] Eric Wallace, Yizhong Wang, Sujian Li, Sameer Singh, and Matt Gardner. Do NLP models know numbers? probing numeracy in embeddings. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 5307–5315, Hong Kong, China, November 2019. Association for Computational Linguistics.
- [113] Zhizhong Li and Derek Hoiem. Learning without forgetting. *IEEE transactions on pattern analysis and machine intelligence*, 40(12):2935–2947, 2017.
- [114] Yizhe Zhang, Siqi Sun, Michel Galley, Yen-Chun Chen, Chris Brockett, Xiang Gao, Jianfeng Gao, Jingjing Liu, and Bill Dolan. Dialogpt: Large-scale generative pre-training for conversational response generation. *arXiv preprint arXiv:1911.00536*, 2019.
- [115] Leonardo FR Ribeiro, Martin Schmitt, Hinrich Schütze, and Iryna Gurevych. Investigating pretrained language models for graph-to-text generation. *arXiv preprint arXiv:2007.08426*, 2020.
- [116] Jeff Donahue, Philipp Krähenbühl, and Trevor Darrell. Adversarial feature learning. *arXiv preprint arXiv:1605.09782*, 2016.
- [117] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pages 1597–1607. PMLR, 2020.
- [118] Ethan Perez, Florian Strub, Harm De Vries, Vincent Dumoulin, and Aaron Courville. Film: Visual reasoning with a general conditioning layer. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32, 2018.
- [119] Ting Chen, Mario Lucic, Neil Houlsby, and Sylvain Gelly. On self modulation for generative adversarial networks. *arXiv preprint arXiv:1810.01365*, 2018.
- [120] Jason Phang, Thibault Févry, and Samuel R Bowman. Sentence encoders on stilts: Supplementary training on intermediate labeled-data tasks. *arXiv preprint arXiv:1811.01088*, 2018.
- [121] Nikita Moghe, Mark Steedman, and Alexandra Birch. Cross-lingual intermediate fine-tuning improves dialogue state tracking. *arXiv preprint arXiv:2109.13620*, 2021.
- [122] Huanru Henry Mao, Bodhisattwa Prasad Majumder, Julian McAuley, and Garrison W Cottrell. Improving neural story generation by targeted common sense grounding. *arXiv preprint arXiv:1908.09451*, 2019.

- [123] Alexander R Fabbri, Simeng Han, Haoyuan Li, Haoran Li, Marjan Ghazvininejad, Shafiq Joty, Dragomir Radev, and Yashar Mehdad. Improving zero and few-shot abstractive summarization with intermediate fine-tuning and data augmentation. *arXiv preprint arXiv:2010.12836*, 2020.
- [124] Raphael Rubino and Eiichiro Sumita. Intermediate self-supervised learning for machine translation quality estimation. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 4355–4360, 2020.
- [125] Mozhdeh Gheini, Xiang Ren, and Jonathan May. Cross-attention is all you need: Adapting pretrained transformers for machine translation. *arXiv preprint arXiv:2104.08771*, 2021.
- [126] Tuan Dinh, Daewon Seo, Zhixu Du, Liang Shang, and Kangwook Lee. Improved input reprogramming for gan conditioning. *arXiv preprint arXiv:2201.02692*, 2022.
- [127] Yen-Chun Chen, Zhe Gan, Yu Cheng, Jingzhou Liu, and Jingjing Liu. Distilling knowledge learned in bert for text generation. *arXiv preprint arXiv:1911.03829*, 2019.
- [128] Sanyuan Chen, Yutai Hou, Yiming Cui, Wanxiang Che, Ting Liu, and Xiangzhan Yu. Recall and learn: Fine-tuning deep pretrained language models with less forgetting. *arXiv preprint arXiv:2004.12651*, 2020.
- [129] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 21(140):1–67, 2020.
- [130] Jiachang Liu, Dinghan Shen, Yizhe Zhang, Bill Dolan, Lawrence Carin, and Weizhu Chen. What makes good in-context examples for gpt-3? *arXiv preprint arXiv:2101.06804*, 2021.
- [131] Teven Le Scao and Alexander M Rush. How many data points is a prompt worth? *arXiv preprint arXiv:2103.08493*, 2021.
- [132] Han Guo, Bowen Tan, Zhengzhong Liu, Eric P Xing, and Zhiting Hu. Text generation with efficient (soft) q-learning. *arXiv preprint arXiv:2106.07704*, 2021.
- [133] Taylor Shin, Yasaman Razeghi, Robert L Logan IV, Eric Wallace, and Sameer Singh. Auto-prompt: Eliciting knowledge from language models with automatically generated prompts. *arXiv preprint arXiv:2010.15980*, 2020.
- [134] Xiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation. *arXiv preprint arXiv:2101.00190*, 2021.
- [135] Brian Lester, Rami Al-Rfou, and Noah Constant. The power of scale for parameter-efficient prompt tuning. *arXiv preprint arXiv:2104.08691*, 2021.
- [136] Xiao Liu, Kaixuan Ji, Yicheng Fu, Zhengxiao Du, Zhilin Yang, and Jie Tang. P-tuning v2: Prompt tuning can be comparable to fine-tuning universally across scales and tasks. *arXiv preprint arXiv:2110.07602*, 2021.
- [137] Vadim Borisov, Tobias Leemann, Kathrin Seßler, Johannes Haug, Martin Pawelczyk, and Gjergji Kasneci. Deep neural networks and tabular data: A survey. *arXiv preprint arXiv:2110.01889*, 2021.
- [138] Inkit Padhi, Yair Schiff, Igor Melnyk, Mattia Rigotti, Youssef Mroueh, Pierre Dognin, Jerret Ross, Ravi Nair, and Erik Altman. Tabular transformers for modeling multivariate time series. In *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 3565–3569. IEEE, 2021.
- [139] Ravid Shwartz-Ziv and Amitai Armon. Tabular data: Deep learning is not all you need. *Information Fusion*, 81:84–90, 2022.
- [140] Hussein Hazimeh, Natalia Ponomareva, Petros Mol, Zhenyu Tan, and Rahul Mazumder. The tree ensemble layer: Differentiability meets conditional computation. In *International Conference on Machine Learning*, pages 4138–4148. PMLR, 2020.

- [141] Sergei Popov, Stanislav Morozov, and Artem Babenko. Neural oblivious decision ensembles for deep learning on tabular data. *arXiv preprint arXiv:1909.06312*, 2019.
- [142] Ira Shavitt and Eran Segal. Regularization learning networks: deep learning for tabular datasets. *Advances in Neural Information Processing Systems*, 31, 2018.
- [143] Arlind Kadra, Marius Lindauer, Frank Hutter, and Josif Grabocka. Regularization is all you need: Simple neural nets can excel on tabular data. *arXiv preprint arXiv:2106.11189*, 2021.
- [144] Xin Huang, Ashish Khetan, Milan Cvitkovic, and Zohar Karnin. Tabtransformer: Tabular data modeling using contextual embeddings. *arXiv preprint arXiv:2012.06678*, 2020.
- [145] Sercan O Arik and Tomas Pfister. Tabnet: Attentive interpretable tabular learning. In *AAAI*, volume 35, pages 6679–6687, 2021.
- [146] Gowthami Somepalli, Micah Goldblum, Avi Schwarzschild, C Bayan Bruss, and Tom Goldstein. Saint: Improved neural networks for tabular data via row attention and contrastive pre-training. *arXiv preprint arXiv:2106.01342*, 2021.
- [147] Tin Kam Ho. Random decision forests. In *Proceedings of 3rd international conference on document analysis and recognition*, volume 1, pages 278–282. IEEE, 1995.
- [148] Tianqi Chen, Tong He, Michael Benesty, Vadim Khotilovich, Yuan Tang, Hyunsu Cho, Kailong Chen, et al. Xgboost: extreme gradient boosting. *R package version 0.4-2*, 1(4):1–4, 2015.
- [149] Scott Reed and Nando De Freitas. Neural programmer-interpreters. *arXiv preprint arXiv:1511.06279*, 2015.
- [150] Nitish Shirish Keskar, Bryan McCann, Lav R Varshney, Caiming Xiong, and Richard Socher. Ctrl: A conditional transformer language model for controllable generation. *arXiv preprint arXiv:1909.05858*, 2019.
- [151] Jack W Rae, Sebastian Borgeaud, Trevor Cai, Katie Millican, Jordan Hoffmann, Francis Song, John Aslanides, Sarah Henderson, Roman Ring, Susannah Young, et al. Scaling language models: Methods, analysis & insights from training gopher. *arXiv preprint arXiv:2112.11446*, 2021.
- [152] Gen Li, Nan Duan, Yuejian Fang, Ming Gong, and Dixin Jiang. Unicoder-vl: A universal encoder for vision and language by cross-modal pre-training. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 11336–11344, 2020.
- [153] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katie Millican, Malcolm Reynolds, et al. Flamingo: a visual language model for few-shot learning. *arXiv preprint arXiv:2204.14198*, 2022.
- [154] 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*, pages 8748–8763. PMLR, 2021.
- [155] Chao Jia, Yafei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc Le, Yun-Hsuan Sung, Zhen Li, and Tom Duerig. Scaling up visual and vision-language representation learning with noisy text supervision. In *International Conference on Machine Learning*, pages 4904–4916. PMLR, 2021.
- [156] Scott Reed, Konrad Zolna, Emilio Parisotto, Sergio Gomez Colmenarejo, Alexander Novikov, Gabriel Barth-Maron, Mai Gimenez, Yury Sulsky, Jackie Kay, Jost Tobias Springenberg, et al. A generalist agent. *arXiv preprint arXiv:2205.06175*, 2022.
- [157] Andrew Jaegle, Sebastian Borgeaud, Jean-Baptiste Alayrac, Carl Doersch, Catalin Ionescu, David Ding, Skanda Koppula, Daniel Zoran, Andrew Brock, Evan Shelhamer, et al. Perceiver io: A general architecture for structured inputs & outputs. *arXiv preprint arXiv:2107.14795*, 2021.

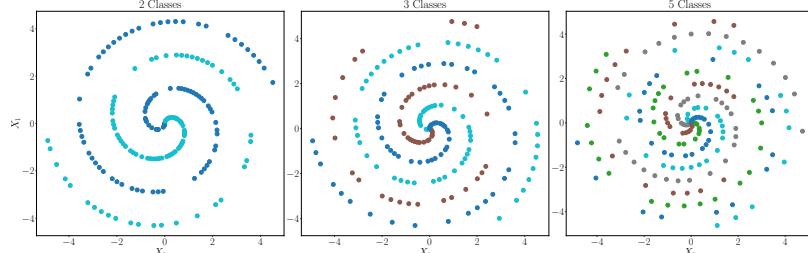
- [158] Matthias Feurer, Aaron Klein, Jost Eggensperger, Katharina Springenberg, Manuel Blum, and Frank Hutter. Efficient and robust automated machine learning. In *Advances in Neural Information Processing Systems 28* (2015), pages 2962–2970, 2015.
- [159] Matthias Feurer, Katharina Eggensperger, Stefan Falkner, Marius Lindauer, and Frank Hutter. Auto-sklearn 2.0: Hands-free automl via meta-learning. *arXiv:2007.04074 [cs.LG]*, 2020.
- [160] Xikun Zhang, Deepak Ramachandran, Ian Tenney, Yanai Elazar, and Dan Roth. Do language embeddings capture scales? *arXiv preprint arXiv:2010.05345*, 2020.
- [161] Aakanksha Naik, Abhilasha Ravichander, Carolyn Rose, and Eduard Hovy. Exploring numeracy in word embeddings. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 3374–3380, 2019.
- [162] Yuanhang Ren and Ye Du. Enhancing the numeracy of word embeddings: A linear algebraic perspective. In *CCF International Conference on Natural Language Processing and Chinese Computing*, pages 170–178. Springer, 2020.
- [163] Guolin Ke, Di He, and Tie-Yan Liu. Rethinking positional encoding in language pre-training. *arXiv preprint arXiv:2006.15595*, 2020.
- [164] Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. Deberta: Decoding-enhanced bert with disentangled attention. *arXiv preprint arXiv:2006.03654*, 2020.
- [165] Benyou Wang, Donghao Zhao, Christina Lioma, Qiuchi Li, Peng Zhang, and Jakob Grue Simonsen. Encoding word order in complex embeddings. *arXiv preprint arXiv:1912.12333*, 2019.
- [166] Zhiheng Huang, Davis Liang, Peng Xu, and Bing Xiang. Improve transformer models with better relative position embeddings. *arXiv preprint arXiv:2009.13658*, 2020.
- [167] Tri Dao, Daniel Y. Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. Flashattention: Fast and memory-efficient exact attention with io-awareness, 2022.
- [168] Emily Sheng, Kai-Wei Chang, Premkumar Natarajan, and Nanyun Peng. The woman worked as a babysitter: On biases in language generation. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3407–3412, Hong Kong, China, November 2019. Association for Computational Linguistics.
- [169] Alekh Agarwal, Alina Beygelzimer, Miroslav Dudík, John Langford, and Hanna Wallach. A reductions approach to fair classification. In *International Conference on Machine Learning*, pages 60–69. PMLR, 2018.
- [170] Yuji Roh, Kangwook Lee, Steven Euijong Whang, and Changho Suh. Fairbatch: Batch selection for model fairness. In *International Conference on Learning Representations*, 2021.
- [171] Michael Gira, Ruisu Zhang, and Kangwook Lee. Debiasing pre-trained language models via efficient fine-tuning. In *Proceedings of the Second Workshop on Language Technology for Equality, Diversity and Inclusion*, pages 59–69, Dublin, Ireland, May 2022. Association for Computational Linguistics.

## Appendix

We first describe details of our datasets in (A.1) and the implementation of all models in (A.2). For the experimental results, we provide the detailed training of LIFT and extended results for several findings presented in the main paper in (B), with more visualizations, score tables, and evaluations.

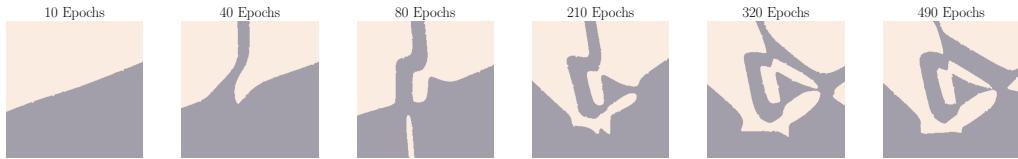
## A Experimental Setup

### A.1 Datasets



**Figure 22:** Rolls dataset of 2, 3, 5 classes. Each dataset contains 300 samples and all classes are balanced.

**Classification datasets.** Table 24 summarizes the datasets used for classification tasks. Besides, we use two additional types of synthetic datasets: **neural-net-based datasets** used for understanding the inductive biases of algorithms (in Sec. 3.2) and **Gaussian pretext datasets** for the two-stage fine-tuning experiments (in Sec. 5.1). The *neural-net-based datasets* are generated as follows. For binary classification, we train a 2-layer neural network with  $\tanh$  activation functions using the *Rolls* dataset shown in the leftmost figure in Fig. 22, and take six snapshots of the decision boundary of the trained neural network shown in Fig. 23; we took snapshots at training epochs 10, 40, 80, 210, 320 and 490. Then, for each snapshot, we define a synthetic dataset (what we call a *neural-net-based dataset*) having labels as the neural network’s prediction for randomly chosen 2000 samples. For 3-class and 5-class classifications, we also use 2000 samples to train a 2-layer neural network using the *Rolls* dataset shown as the second and the third figure in Fig. 22. The decision boundaries of networks trained on more epochs are visually more complex. Hence, the corresponding classification tasks are becoming more complicated, from the left column to the right column in Fig. 23. In the manuscript, we tested on three out of six datasets, obtained by snapshots at 10, 80, and 490 epochs, respectively. Given a target dataset of  $n$  classes and  $d$  features, *Gaussian pretext datasets* is generated as follows: using *scikit-learn*<sup>3</sup>, we randomly generate datasets of  $n$  clusters, where each cluster has 100 normally distributed samples in the  $d$ -dimensional space.



**Figure 23:** Neural-net-based datasets. Given Rolls dataset in Fig. 22, we train a 2-layer neural network for 10, 40, 80, 210, 320, 490 epochs, and get six decision boundaries at each column. We define six neural-net-based datasets from here: each decision boundary is used as a labeling function of each neural-net-based dataset. In the main manuscript, we used three out of six datasets, obtained by snapshots at 10, 80, and 490 epochs.

**Regression datasets.** We test LIFT/GPT on regression problems for both synthetic/real datasets. **(Synthetic datasets)** To assess the regression performance of LIFT in different datasets, we generate synthetic datasets based on six different functions types, including smooth functions, non-smooth functions, and non-continuous functions:

1. Linear functions  $y = f(\mathbf{x}) = \mathbf{x}^\top \mathbf{1}/p$

<sup>3</sup>[https://scikit-learn.org/stable/modules/generated/sklearn.datasets.make\\_classification.html](https://scikit-learn.org/stable/modules/generated/sklearn.datasets.make_classification.html)

Data Type	Dataset	ID	Abbreviation	No. Features	No. Classes	No. Instances	Note
Synthetic	9Gaussians	1	-	2	9	2000	-
	Blobs	2	-	2	4	2000	-
	Circle	3	-	2	2	2000	non-linear boundary
	TwoCircles	6	-	2	2	2000	non-linear boundary
	Moons	4	-	2	2	2000	-
Tabular (OpenML)	wholesale-customers	1511	Customers	8	2	440	Imbalance
	pollution	882	Pollution	15	2	60	1 symbolic feature
	spambase	44	Spambase	57	2	4601	1 symbolic feature
	hill-valley	1479	Hill-Valley	100	2	1212	1 symbolic feature
	tae	48	TAE	5	3	151	Categorical data
	cmc	23	CMC	9	3	1473	Meaningful feature Names
	wine	187	Wine	13	3	178	Integral features
	vehicle	54	Vehicle	18	4	846	Meaningful feature Names
	LED-display-domain-7digit	40496	LED	7	10	500	1 symbolic feature
	optdigits	28	OPT	64	10	5620	1 symbolic feature
	mfeat-factors	12	Mfeat	216	10	2000	1 symbolic feature
	pollen	871	Pollen	5	2	3848	-
	climate-model-simulation-crashes	1467	Climate	20	2	540	-
	one-hundred-plants-margin	1491	Margin	64	100	1600	1 symbolic feature
	one-hundred-plants-shape	1492	Shape	64	100	1600	1 symbolic feature
	one-hundred-plants-texture	1493	Texture	64	100	1599	1 symbolic feature
	breast-cancer	13	Breast	9	2	286	-
	iris	61	Iris	4	3	150	-
	visualizing_hamster	893	Hamster	5	2	73	-
	PizzaCutter3	1444	Pizza	37	2	1043	-
Vision	MNIST	-	-	784	10	70k	-
	Permuted MNIST	-	P-MNIST	784	10	70k	-
	Fashion MNIST	-	FMNIST	784	10	70k	-
	Permuted Fashion MNIST	-	P-FMNIST	784	10	70k	-

**Table 24: Classification datasets.** We have three non-language types of data: synthetic data, real tabular data and vision data. We use five synthetic datasets. For the real tabular data, we select datasets from OpenML with a wide range of number of features, types of features, number of classes, and number of training samples. We use MNIST, Fashion MNIST, and their permuted variants for the vision datasets.

2. Quadratic functions  $y = f(\mathbf{x}) = \mathbf{x}^T \mathbf{I} \mathbf{x} / p$ , where  $\mathbf{I}$  is the identity matrix
3. Continuous exponential function  $y = f(\mathbf{x}) = \sum_{i=1}^p e^{0.2x_i} / p$
4. Cosine functions  $y = f(\mathbf{x}) = \sum_{i=1}^p \cos(0.5\pi x_i) / p$
5. (Non-smooth)  $\ell_1$ -norm function  $y = f(\mathbf{x}) = \|\mathbf{x}\|_1 / p$
6. (Non-continuous) Piecewise linear function

$$f(\mathbf{x}) = \frac{1}{p} \sum_{i=1}^p \tilde{f}(x_i) := \begin{cases} x_i - 1 & -10 \leq x_i < -3, \\ 0 & -3 \leq x_i < 3, \\ x_i + 1 & 3 \leq x_i \leq 10. \end{cases}$$

We let  $x_i \sim \text{Unif}(-10, 10)$  for each coordinate  $i$ , and the noise level  $\sigma = 0.1$  by default. We normalize all the functions above for fair comparison among different functions so that their  $y \in [-9, 9]$  when  $\mathbf{x} \in [-10, 10]^p$ . In particular, to assess whether LIFT is better at dealing with positive numbers or integers, we generate additional datasets by further manipulating the  $(\mathbf{x}, \mathbf{y})$  distribution of linear and piecewise functions. For datasets with real numbers, we generate the 1D dataset  $\mathbf{x} \sim \text{Unif}(-150, -150)$ . For datasets with only positive numbers, we generate the dataset  $\mathbf{x}_i \sim \text{Unif}(0, 300)$ . To generate the datasets with all integer prompts, we round down all the features to integers. For visualization, in addition to the training, validation, and test datasets, we generated grid datasets. Unless otherwise stated, we generate uniformly spaced 200 samples for 1D visualizations and 2,500 samples for 2D visualizations, with each coordinate  $x_i \in [-10, 10]$  for  $i = 1, \dots, p$ . To visualize the extrapolation performance, we let the  $x_i \in [-15, 15]$ .

**(Real datasets)** We consider four different types of real datasets: Medical Insurance dataset [41] with 1,338 samples and 6 features, Combined Cycle Power Plant (CCPP) dataset [42] with 9,568 samples and 4 features, Servo [43] dataset with 167 samples and 4 features, and Student [44] dataset with 649 samples and 33 features. In particular, the Medical Insurance dataset and Student dataset contain feature names that can be interpreted using common knowledge (see feature lists in Table 33), while CCPP and Servo do not.

## A.2 LIFT and Baseline Implementation

This section provides details of our models and implementation. We describe our pretrained language models and the baseline implementations (A.2.1), the computing resources used for running experiments (A.2.2), and how to fine-tune and select the hyperparameters (A.2.3).

### A.2.1 Pretrained Language Models and Baselines

**Pretrained language models.** Our main results are with two pretrained language models: GPT-J [31] and GPT-3 [19]. We mainly focus on GPT-J for the reproducibility purpose and provide additional results on GPT-3 as reference. For GPT-J, we use a quantized version<sup>4</sup> of 6 billion parameters with 8-bit weights. For GPT-3, we use the GPT-3 OpenAI API<sup>5</sup> and employ the Ada version by default. In Sec. 3.6, we compare two previous mentioned models with three bigger versions of GPT-3 (*Baggage*, *Curie*, *Davinci*). The largest one is *Davinci*-GPT-3 containing approximately 175 billions of parameters.

Since GPTs can output any language token, the output might not be appropriate for the desired task. For example, GPT might output non-numerical words for regression task approximating a function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ . In such a case, we categorize this as *invalid output*.

**Baselines** For XGBoost, we use the open-source XGBoost module<sup>6</sup>. For other baselines, we implement them using scikit-learn module<sup>7</sup>.

### A.2.2 Computing Resources

For experiments on GPT-J, we use GPU computing from two 24GB-RTX3090 GPUs and AWS EC2 instances<sup>8</sup> (p3.8xlarge, p3.2xlarge). For other models, we run experiments on CPU instances.

### A.2.3 Hyperparameter Selection

In fine-tuning GPT-J, we use Adam-8bit optimizer implementation<sup>9</sup> with weight decay of 0.01 and 6 warm-up steps. The learning rate is chosen from  $1e-4$  and  $2e-4$  for the synthetic/OpenML datasets, and  $1e-5$  for the vision datasets. We use a linear learning scheduler for the optimizer. For classification, the batch size depends on the number of features of the datasets. We set the batch size to be 128, 32, 16, and 2 for datasets with the number of features being no greater than 2, between 2 and 6, between 6 and 20, and greater than 20, respectively. For regression, we set batch size as 4 by default and reduce it to 1 to avoid the memory issue when the number of features increases. For GPT-3, we use the API provided by OpenAI to perform black-box GPT-3 fine-tuning with the default setting. Our implementation is with PyTorch framework<sup>10</sup>.

We perform hyperparameter selection based on validation results for all methods for a fair comparison. The hyperparameter tuning scheme for all methods is detailed as follows:

#### Classification methods

- LIFT/GPT-J: number of epochs  $\in \{5, 10, 15\}$
- LIFT/GPT-3: learning rate multiplier  $\in \{0.05, 0.1, 0.2\}$
- Random Forest (RF): maximum depth  $\in \{3, 5, 10\}$ , minimum number of samples required to split an internal node  $\in \{2, 5, 10\}$
- Decision Tree (DT): maximum depth of the tree  $\in \{3, 5, 20\}$  and criterion  $\in \{\text{Gini impurity}, \text{Shannon information gain}\}$

<sup>4</sup><https://huggingface.co/hivemind/gpt-j-6B-8bit>

<sup>5</sup><https://openai.com/api/>

<sup>6</sup><https://xgboost.readthedocs.io/>

<sup>7</sup><https://scikit-learn.org/>

<sup>8</sup><https://aws.amazon.com/ec2/>

<sup>9</sup><https://huggingface.co/hivemind/gpt-j-6B-8bit>

<sup>10</sup><https://pytorch.org/>

- Support Vector Machine (*SVM*): kernel  $\in \{\text{linear kernel, radial basis function}\}$ , regularization parameter  $\in \{1, 10, 100\}$
- Multilayer Perceptron (*MLP*): initial learning rate  $\in \{0.001, 0.01, 0.1\}$
- Logistic regression (*LogReg*): inverse of regularization strength  $\in \{1, 10, 100\}$
- K-Nearest Neighbor (*KNN*): number of neighbors to use  $\in \{1, 3, 5\}$ , power parameter for the Minkowski metric  $\in \{1, 2\}$
- XGBoost (*XG*): maximum depth  $\in \{3, 5, 10\}$

### Regression methods

- LIFT/GPT-J: number of epochs  $\in \{2, 6, 10\}$
- LIFT/GPT-3: learning rate multiplier  $\in \{0.05, 0.1, 0.2\}$
- Polynomial Regression (*PR*): no hyperparameter selection but fixed the degree at 3 since higher-order polynomial regression introduces out-of-memory error especially for high-dimensional datasets
- K-Nearest Neighbor (*KNN*): number of neighbors  $\in \{2, 5, 8\}$
- Kernel Regression (*KR*): Gamma parameter of Radial Basis Kernel  $\in \{0.01, 0.1, 1\}$
- Multilayer Perceptron (*MLP*): initial learning rate  $\in \{0.0001, 0.001, 0.01\}$
- Gradient Boosting Tree (*GBT*): learning rate  $\in \{0.001, 0.01, 0.1\}$
- Random Forest (*RF*): maximum depth  $\in \{4, 6\}$
- Gaussian Process (*GP*): the number of optimizer restarts used to find parameters of the kernel that maximize the log marginal likelihood  $\in \{5, 10\}$ .

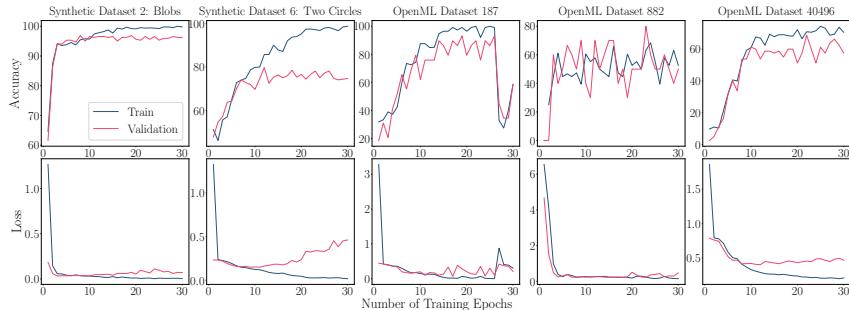
Note that we perform model selection based on validation RAE, instead of validation loss.

## B Detailed and Extended Experimental Results

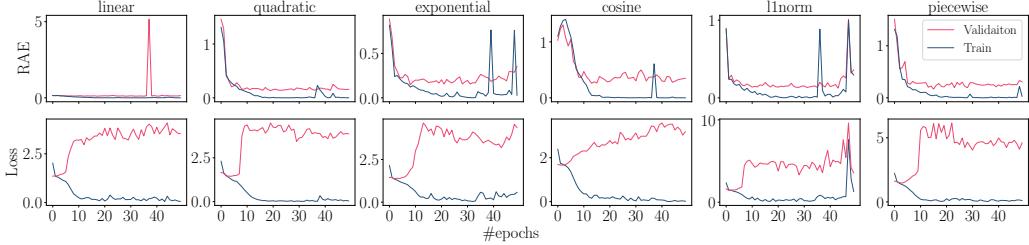
We provide the extended version of the experimental results presented in the main paper.

### B.1 LIFT’s Training

We report the learning curves of LIFT in terms of LM-loss and accuracies/RAE for several classification and regression tasks. We observe the decrease of training loss over the tasks and datasets. We select the best models based on the validation criteria (accuracy for classification and RAE for regression) on the validation sets. Fig. 25 visualize the accuracy and loss of LIFT/GPT-J in the training and validation process for classification tasks. For regression task, Fig. 26 shows that the decrease in RAE does not necessarily imply the decrease in loss. Furthermore, we observe that LIFT only requires a few epochs to achieve good performance.



**Figure 25: Learning curves of LIFT/GPT-J on several synthetic/OpenML classification datasets.** We plot the accuracy (top row) and the loss (bottom row) of LIFT varying the number of training epochs.



**Figure 26: Regression RAE and Loss curves of LIFT/GPT-J on the synthetic regression datasets.** We observe that LIFT/GPT-J only requires a few epochs to achieve good performance.

## B.2 Extended Results

### B.2.1 Classification and Regression Evaluations on All Baseline Methods

We provide the full results of classification and regression performances with all baselines. Table 27 presents the classification performance with other considered baselines, including KNN, MLP, and Random Forest.

Table 28 provides the regression evaluation with all regression baselines on synthetic datasets, and Table 29 provides the results for real datasets. Since experiments with LIFT/GPT-J are conducted on AWS and local server and due to this limitation of memory resources, we fail to run experiments of LIFT/GPT-J on high-dimensional datasets. Therefore, for 50D and 120D synthetic datasets, only results of LIFT/GPT-3 are reported.

Type	Dataset (ID)	p/c	ODC	LogReg	KNN	DT	MLP	RBF-SVM	RF	XG	LIFT/GPT-J	LIFT/GPT-3
Synthetic	circles (3)	2 / 2	50.00	48.58±1.94	81.25±0.20	77.42±0.24	82.00±0.54	<b>83.08±0.59</b>	82.42±1.33	81.42±0.31	79.95±1.53	81.17±0.42
	two circles (6)	2 / 2	50.00	49.83±4.18	<b>81.83±0.62</b>	75.50±0.20	68.42±3.86	80.00±0.54	76.08±0.59	79.25±0.35	75.92±1.65	81.42±0.82
	blobs (2)	2 / 4	25.00	96.75±0.00	95.50±0.20	96.08±0.82	96.58±0.42	96.75±0.00	<b>97.17±0.24</b>	96.17±0.12	96.17±0.59	96.67±0.24
	moons (4)	2 / 4	50.00	88.58±0.12	<b>100.00±0.00</b>	99.25±0.41	98.75±1.08	<b>100.00±0.00</b>	99.75±0.00	99.83±0.12	99.58±0.42	<b>100.00±0.00</b>
	9Clusters (1)	2 / 9	11.25	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	100.00±0.00	<b>100.00±0.00</b>	<b>100.00±0.00</b>	<b>100.00±0.00</b>	<b>100.00±0.00</b>
	Customers (1511)	9 / 2	68.18	87.12±0.54	<b>88.64±0.00</b>	85.98±0.53	86.36±1.86	86.36±0.00	85.23±0.00	85.23±0.00	85.23±1.61	84.85±1.42
Tabular (OpenML)	Pollution (882)	16 / 2	50.00	58.33±11.79	<b>66.67±6.81</b>	<b>77.78±3.93</b>	66.67±0.00	58.33±6.81	<b>77.78±3.93</b>	63.89±7.86	63.89±3.93	63.89±7.86
	Spambase (44)	58 / 2	60.59	93.27±0.00	90.77±0.00	90.77±0.14	<b>94.35±0.00</b>	93.70±0.00	95.01±0.00	<b>95.87±0.00</b>	94.03±0.54	94.90±0.36
	Hill-Valley (1479)	101 / 2	49.79	77.78±0.00	56.38±0.00	56.38±0.89	50.21±0.00	68.72±0.00	51.44±0.00	59.26±0.00	<b>100.00±0.20</b>	<b>99.73±0.19</b>
	TAE (48)	6 / 3	35.48	45.16±4.56	60.22±4.00	65.59±5.49	54.84±2.63	53.76±6.63	<b>67.74±7.90</b>	66.67±8.05	61.29±6.97	65.59±6.63
	CMC (23)	10 / 3	42.71	49.49±0.83	50.85±1.91	56.72±0.32	57.29±0.73	56.50±0.97	53.45±1.05	52.43±0.42	49.83±0.28	<b>57.74±0.89</b>
	Wine (187)	14 / 3	38.89	<b>100.00±0.00</b>	96.20±1.31	93.52±2.62	98.15±2.52	<b>100.00±0.00</b>	<b>100.00±0.00</b>	97.22±0.00	93.52±1.31	92.59±1.31
	Vehicle (34)	19 / 4	25.88	80.39±1.00	69.61±0.74	63.92±2.7	79.21±0.28	<b>81.18±0.48</b>	75.88±1.27	73.14±0.28	64.31±2.37	70.20±2.73
	LED (40496)	8 / 10	11.00	68.67±0.94	63.67±6.13	66.33±2.87	<b>72.00±0.82</b>	68.00±0.82	64.33±0.94	66.00±0.82	65.33±0.47	69.33±2.05
	OPT (28)	65 / 10	10.14	96.53±0.22	96.92±0.16	89.8±1.09	97.36±0.27	97.95±0.00	97.69±0.14	97.48±0.17	98.22±0.11	<b>98.99±0.30</b>
	Mfeat (12)	217 / 10	10.00	97.67±0.12	97.67±0.3	87.67±1.05	96.5±0.35	<b>98.83±0.24</b>	97.75±0.35	96.75±0.00	94.17±1.75	93.08±0.24
	Margin (1491)	65 / 100	0.94	81.35±0.15	77.60±0.97	43.86±0.21	77.71±1.91	<b>81.98±0.30</b>	77.71±1.98	70.21±0.29	50.23±1.33	59.37±0.92
	Texture (1493)	65 / 100	0.94	81.67±0.97	80.62±0.76	46.88±1.93	76.88±2.44	<b>83.44±0.89</b>	73.12±0.76	70.73±1.41	50.32±2.18	67.50±1.42
Images	MNIST		11.35	91.95±0.69	96.71±0.11	87.42±0.64	97.30±0.16	97.70±0.97	94.91±0.18	97.69±0.04	97.01±1.15	<b>98.15±0.67</b>
	P-MNIST	784 / 10	11.35	92.58±0.04	96.74±0.08	87.87±0.69	97.39±0.14	<b>98.06±0.31</b>	94.59±0.18	97.62±0.09	95.80±0.07	96.25±0.35
	FMNIST		10.00	85.59±0.09	85.59±0.03	80.52±0.40	88.86±0.02	<b>90.59±0.02</b>	85.25±0.13	90.19±0.04	85.10 ± 0.19	90.18 ± 0.12
	P-FMNIST		10.00	84.95±0.84	85.15±0.61	79.91±0.93	88.86±0.61	88.04±0.69	84.93±0.59	<b>89.93±0.14</b>	82.25±0.27	88.92±0.71

**Table 27: Accuracies (↑) on various classification datasets.** We evaluate LIFT/GPTs on different classification datasets: 2D synthetic datasets, tabular datasets in OpenML [36], and image datasets, varying number of features ( $p$ ) and number of classes ( $c$ ). Overall, LIFT/GPTs perform comparably well across all tasks. LIFT/GPTs can be adapted well to non-linear datasets (circles, two circles), beyond the capacity of logistic regression. On the OpenML data, LIFT/GPTs achieves competitive performances with the best methods, such as XGBoost or RBF-SVM. The performance of LIFT/GPTs degrades as the number of classes is large, e.g., when the number of classes  $c=100$ . On the vision data, LIFT/GPTs perform relatively well, achieving highly competitive accuracies on both MNIST and Fashion MNIST. We note that the classes of MNIST are not fully balanced. Thus ODC gets 11.35% instead of 10% as ODC returns the optimal class learned from the training dataset.

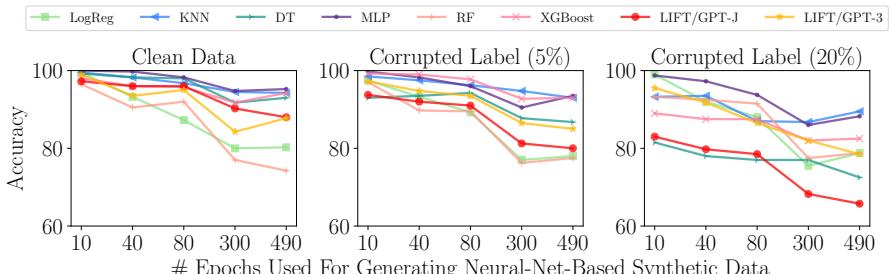
Dataset \ Method	PR	KR	KNN	MLP	GBT	RF	GP	LIFT/GPT-J	LIFT/GPT-3
Linear	$p = 1$ $0.01 \pm 0.0$	$0.05 \pm 0.0$	$0.04 \pm 0.0$	$0.03 \pm 0.0$	$0.05 \pm 0.0$	$0.04 \pm 0.0$	$0.01 \pm 0.0$	$0.08 \pm 0.0$	$0.06 \pm 0.0$
	$p = 2$ $0.03 \pm 0.0$	$0.09 \pm 0.0$	$0.12 \pm 0.0$	$0.04 \pm 0.0$	$0.12 \pm 0.0$	$0.12 \pm 0.0$	$0.01 \pm 0.0$	$0.12 \pm 0.0$	$0.19 \pm 0.0$
	$p = 50$ $0.71 \pm 0.0$	$1.02 \pm 0.0$	$0.78 \pm 0.0$	$1.85 \pm 0.1$	$0.97 \pm 0.0$	$0.87 \pm 0.0$	$0.13 \pm 0.0$	-	$1.18 \pm 0.2$
	$p = 100$ $0.95 \pm 0.0$	$1.02 \pm 0.0$	$0.88 \pm 0.0$	$3.02 \pm 0.0$	$0.99 \pm 0.0$	$0.94 \pm 0.0$	$0.64 \pm 0.0$	-	$2.14 \pm 0.5$
Quadratic	$p = 1$ $0.01 \pm 0.0$	$0.05 \pm 0.0$	$0.05 \pm 0.0$	$0.03 \pm 0.0$	$0.06 \pm 0.0$	$0.05 \pm 0.0$	$0.01 \pm 0.0$	$0.11 \pm 0.0$	$0.13 \pm 0.0$
	$p = 2$ $0.03 \pm 0.0$	$0.16 \pm 0.0$	$0.17 \pm 0.0$	$0.06 \pm 0.0$	$0.15 \pm 0.0$	$0.25 \pm 0.0$	$0.02 \pm 0.0$	$0.28 \pm 0.1$	$0.22 \pm 0.0$
	$p = 50$ $1.12 \pm 0.0$	$5.19 \pm 0.0$	$1.33 \pm 0.0$	$2.28 \pm 0.0$	$0.98 \pm 0.0$	$0.96 \pm 0.0$	$0.69 \pm 0.0$	-	$0.99 \pm 0.2$
	$p = 100$ $1.02 \pm 0.0$	$7.30 \pm 0.0$	$1.29 \pm 0.0$	$2.89 \pm 0.0$	$1.01 \pm 0.0$	$0.98 \pm 0.0$	$0.89 \pm 0.0$	-	$1.06 \pm 0.1$
Exponential	$p = 1$ $0.04 \pm 0.0$	$0.07 \pm 0.0$	$0.05 \pm 0.0$	$0.02 \pm 0.0$	$0.05 \pm 0.0$	$0.04 \pm 0.0$	$0.01 \pm 0.0$	$0.11 \pm 0.0$	$0.09 \pm 0.0$
	$p = 2$ $0.04 \pm 0.0$	$0.15 \pm 0.0$	$0.13 \pm 0.0$	$0.07 \pm 0.0$	$0.09 \pm 0.0$	$0.11 \pm 0.0$	$0.04 \pm 0.0$	$0.19 \pm 0.0$	$0.20 \pm 0.0$
	$p = 50$ $0.94 \pm 0.0$	$10.23 \pm 0.0$	$1.04 \pm 0.0$	$3.18 \pm 0.2$	$1.05 \pm 0.0$	$0.96 \pm 0.0$	$0.53 \pm 0.0$	-	$1.15 \pm 0.0$
	$p = 100$ $0.96 \pm 0.0$	$14.12 \pm 0.0$	$1.03 \pm 0.0$	$4.14 \pm 0.0$	$0.97 \pm 0.0$	$0.93 \pm 0.0$	$0.79 \pm 0.0$	-	$1.03 \pm 0.0$
Cosine	$p = 1$ $1.05 \pm 0.0$	$0.12 \pm 0.0$	$0.14 \pm 0.0$	$0.38 \pm 0.1$	$0.15 \pm 0.0$	$0.35 \pm 0.0$	$0.04 \pm 0.0$	$0.38 \pm 0.1$	$0.44 \pm 0.1$
	$p = 2$ $1.04 \pm 0.0$	$0.74 \pm 0.0$	$0.83 \pm 0.1$	$1.06 \pm 0.0$	$0.41 \pm 0.0$	$0.80 \pm 0.0$	$0.31 \pm 0.0$	$0.82 \pm 0.2$	$0.65 \pm 0.1$
	$p = 50$ $1.01 \pm 0.0$	$1.01 \pm 0.0$	$1.00 \pm 0.0$	$1.59 \pm 0.0$	$1.00 \pm 0.0$	$0.99 \pm 0.0$	$1.01 \pm 0.0$	-	$1.25 \pm 0.1$
	$p = 100$ $1.02 \pm 0.0$	$1.00 \pm 0.0$	$1.09 \pm 0.0$	$2.43 \pm 0.1$	$1.04 \pm 0.0$	$1.06 \pm 0.0$	$1.00 \pm 0.0$	-	$1.20 \pm 0.3$
L1norm	$p = 1$ $0.23 \pm 0.0$	$0.06 \pm 0.0$	$0.05 \pm 0.0$	$0.03 \pm 0.0$	$0.06 \pm 0.0$	$0.06 \pm 0.0$	$0.03 \pm 0.0$	$0.10 \pm 0.0$	$0.09 \pm 0.0$
	$p = 2$ $0.24 \pm 0.0$	$0.17 \pm 0.0$	$0.19 \pm 0.0$	$0.06 \pm 0.0$	$0.15 \pm 0.0$	$0.29 \pm 0.0$	$0.07 \pm 0.0$	$0.24 \pm 0.0$	$0.20 \pm 0.0$
	$p = 50$ $1.09 \pm 0.0$	$1.00 \pm 0.0$	$1.28 \pm 0.0$	$1.97 \pm 0.1$	$0.98 \pm 0.0$	$0.94 \pm 0.0$	$0.96 \pm 0.0$	-	$1.12 \pm 0.1$
	$p = 100$ $1.01 \pm 0.0$	$1.01 \pm 0.0$	$1.22 \pm 0.0$	$2.80 \pm 0.1$	$1.03 \pm 0.0$	$1.01 \pm 0.0$	$0.99 \pm 0.0$	-	$1.27 \pm 0.2$
Piecewise	$p = 1$ $0.45 \pm 0.0$	$0.17 \pm 0.0$	$0.08 \pm 0.0$	$0.08 \pm 0.0$	$0.06 \pm 0.0$	$0.07 \pm 0.0$	$0.10 \pm 0.0$	$0.15 \pm 0.0$	$0.17 \pm 0.0$
	$p = 2$ $0.39 \pm 0.0$	$0.34 \pm 0.0$	$0.33 \pm 0.0$	$0.20 \pm 0.0$	$0.19 \pm 0.0$	$0.38 \pm 0.0$	$0.29 \pm 0.0$	$0.40 \pm 0.1$	$0.40 \pm 0.1$
	$p = 50$ $0.93 \pm 0.0$	$1.00 \pm 0.0$	$0.97 \pm 0.0$	$2.11 \pm 0.0$	$1.00 \pm 0.0$	$0.94 \pm 0.0$	$0.93 \pm 0.0$	-	$1.35 \pm 0.1$
	$p = 100$ $1.01 \pm 0.0$	$1.00 \pm 0.0$	$1.08 \pm 0.0$	$4.20 \pm 0.1$	$1.02 \pm 0.0$	$1.01 \pm 0.0$	$1.01 \pm 0.0$	-	$1.11 \pm 0.0$

**Table 28: Comparison of regression methods in approximating various functions.** The regression performance is measured by RAE( $\downarrow$ ), and we tested on six functions with various  $p$ , the number of features. LIFT can approximate different types of functions in low-dimensional cases ( $p = 1, 2$ ), although it fails to achieve performance comparable to that of strong baselines. We observed that LIFT fails to achieve satisfying regression performance in high-dimensional cases ( $p = 50, 100$ ). Due to resource limitations, we have not obtained results of LIFT/GPT-J on high-dimensional datasets.

Dataset \ Method	PR	KR	KNN	MLP	GBT	RF	GP	LIFT/GPT-J	LIFT/GPT-3
ccpp	$0.22 \pm 0.00$	$21.60 \pm 0.00$	$0.45 \pm 0.00$	$0.30 \pm 0.00$	$0.17 \pm 0.00$	$0.21 \pm 0.00$	$0.69 \pm 0.00$	$0.24 \pm 0.01$	$0.18 \pm 0.01$
servo	$0.92 \pm 0.00$	$0.95 \pm 0.00$	$0.86 \pm 0.00$	$0.82 \pm 0.00$	$0.25 \pm 0.00$	$0.25 \pm 0.00$	$1.03 \pm 0.00$	$1.17 \pm 0.16$	$0.29 \pm 0.02$
insurance	$0.48 \pm 0.00$	$1.48 \pm 0.00$	$1.03 \pm 0.00$	$0.44 \pm 0.00$	$0.25 \pm 0.00$	$0.26 \pm 0.00$	$1.30 \pm 0.00$	$0.53 \pm 0.11$	$0.14 \pm 0.05$
student	$0.47 \pm 0.00$	$1.56 \pm 0.00$	$0.66 \pm 0.00$	$0.37 \pm 0.00$	$0.39 \pm 0.00$	$0.36 \pm 0.00$	$0.45 \pm 0.00$	$0.36 \pm 0.02$	$0.27 \pm 0.01$

**Table 29: Comparison of regression methods in real datasets.** The regression performance is measured by RAE( $\downarrow$ ). We observe that LIFT/GPT-3 achieves the top 2 regression performance among all the real datasets.

### B.2.2 Quantitative Classification Evaluations on Neural-Net-Based Synthetic Datasets

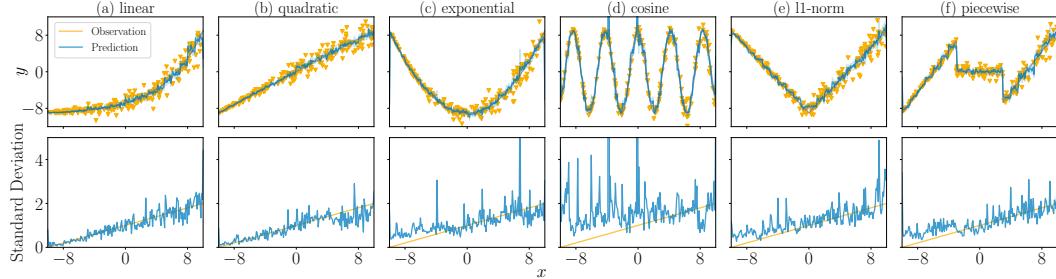


**Figure 30: How accuracy ( $\uparrow$ ) changes as the target classification problem becomes more complex, i.e., the ground-truth decision boundary becomes more complex.** The x-axis shows the number of epochs we used to train the neural network on the Rölls dataset. Note that the network becomes more complex as the number of epochs increases. Thus the classification problem also gets challenging. We measure the performances on three cases: (left) clean data, (middle) label-corrupted data with corruption probability 5% and (right) 20%.

In Sec. 3.2, we assess how well LIFT/GPTs adapt to different shapes of decision boundaries on neural-net-based synthetic datasets. Here, we provide the testing accuracies of all models on these datasets. For binary-class datasets, Fig. 30 shows the accuracies on three binary-class datasets when the difficulty of classification tasks varies by using different network checkpoints at different epochs. As the difficulty increases or the level of corruption increases, all methods tend to decrease classification accuracy. We provide three settings of training data: clean data, 5% label corruption, and 20% label corruption. We observe that LIFT/GPT-3 outperforms logistic regression and decision tree,

especially when the label corruption is up to 20%. However, LIFT/GPT-J is performing worse than other baselines in the data corruption scenarios. For 3-class and 5-class datasets, both LIFT/GPT-3 and LIFT/GPT-3 achieve approximately 90%, while the best baselines (MLP and XGBoost) obtains approximately 92% and 91% for the 3-class and 5-class data, respectively.

### B.2.3 Estimating Predictive Uncertainty with LIFT/GPT-3



**Figure 31: Visualization of LIFT/GPT-3 predictions under varying noise levels.** The predictions are made on grid datasets consisting of 103 evenly-spaced samples in  $[-10, 10]$ . The standard deviation of LIFT/GPT-3 predictions are computed based 20 repeated random predictions. We observe that the standard deviations of predictions from LIFT/GPT-3 aligns well with that of noisy training samples (observations), implying that LIFT/GPT-3 can quantify the predictive uncertainty.

Continuing the discussion in Sec. 4.2, Fig. 31 indicates that LIFT/GPT-3 also shares similar behaviors to Bayesian estimators in quantifying the uncertainty of the predictions.

### B.2.4 Does LIFT Benefit from Incorporating Feature Names?

**Classification.** We provide additional evaluation of LIFT/GPT-J on three datasets used in the main paper. Table 32 presents our result with the same settings in the main paper. We can see that correctly using feature names helps improve the performance of LIFT from the models without feature names or the models with randomly shuffled feature names. This finding is consistent with the finding in the main papers on the usefulness of incorporating the feature names.

Dataset (ID)	ODC	XGBoost	LIFT/GPT-J				
			W/o Names	Shuffled-Names I	Shuffled-Names II	Correct-Names I	Correct-Names II
CMC (23)	42.71	$52.43 \pm 0.42$	$49.49 \pm 0.56$	<b><math>51.30 \pm 1.05</math></b>	<b><math>51.30 \pm 2.51</math></b>	$48.82 \pm 3.12$	$50.39 \pm 1.05$
TAE (48)	35.48	$66.67 \pm 8.05$	$60.22 \pm 4.02$	$63.44 \pm 6.08$	$58.06 \pm 7.90$	$60.21 \pm 10.64$	<b><math>65.59 \pm 8.47</math></b>
Vehicle (54)	25.88	$73.14 \pm 0.28$	$64.31 \pm 2.37$	$66.87 \pm 1.54$	$65.49 \pm 1.69$	<b><math>69.02 \pm 3.67^*</math></b>	

**Table 32: The effect of using feature names on LIFT/GPT-J.** We compare classification accuracy ( $\uparrow$ ) of LIFT/GPT-J when feature names provided in the target dataset *are and are not* incorporated to the prompts.

**Regression.** To investigate whether incorporating feature names in prompts improves the regression performance of LIFT, similar to the datasets selection process of classification tasks, we evaluate the effect of feature names on the datasets Insurance and Student, whose tasks can be helped by common knowledge. To be more specific, while the task of Insurance dataset is to predict the insurance costs, the key features of Insurance dataset are age, body mass index, and smoke or not, which are intuitively closely related to the task. For the Student dataset, the task is to predict students’ grades based on their weekly study time, previous grades, etc. Therefore, the features and task of Student are also highly correlated. Table 33 presents our evaluation of regression tasks. We find that fine-tuning with feature names does not necessarily help regression tasks.

### B.2.5 Robustness to Label Corruption

Table 34 and Table 35 extend the results reported in Fig. 11. These additional datasets follow a similar trend to what was discussed in Sec. 3.4.

Dataset	Frac.	RF	LIFT/GPT-3				
			W/O Names	Shuffled-Names I	Shuffled-Names II	Correct-Names I	Correct-Names II
insurance	0.2	<b>0.31 ± 0.00</b>	0.89 ± 0.03	0.76 ± 0.11	0.59 ± 0.09	0.59 ± 0.11	0.89 ± 0.03
	0.4	0.26 ± 0.00	0.42 ± 0.15	0.30 ± 0.02	<b>0.20 ± 0.03</b>	0.35 ± 0.10	0.21 ± 0.01
	0.6	0.26 ± 0.00	0.30 ± 0.10	0.24 ± 0.03	<b>0.19 ± 0.02</b>	0.30 ± 0.12	0.22 ± 0.08
	0.8	0.27 ± 0.00	0.31 ± 0.07	0.19 ± 0.04	0.18 ± 0.03	0.14 ± 0.01	<b>0.11 ± 0.02</b>
	1.0	0.26 ± 0.00	0.14 ± 0.05	0.17 ± 0.03	0.19 ± 0.01	0.17 ± 0.04	<b>0.10 ± 0.03</b>
student	0.2	0.40 ± 0.00	0.32 ± 0.01	0.32 ± 0.01	0.34 ± 0.02	<b>0.31 ± 0.01</b>	<b>0.31 ± 0.01</b>
	0.4	0.36 ± 0.00	0.32 ± 0.02	0.31 ± 0.01	<b>0.30 ± 0.00</b>	0.32 ± 0.01	0.35 ± 0.01
	0.6	0.36 ± 0.00	0.31 ± 0.01	0.31 ± 0.01	0.31 ± 0.01	0.31 ± 0.01	<b>0.30 ± 0.00</b>
	0.8	0.38 ± 0.00	0.28 ± 0.01	<b>0.27 ± 0.01</b>	0.29 ± 0.02	0.28 ± 0.01	0.28 ± 0.00
	1.0	0.35 ± 0.00	<b>0.27 ± 0.01</b>	0.28 ± 0.01	0.28 ± 0.01	0.28 ± 0.01	0.35 ± 0.02

**Table 33: Investigating if incorporating feature names to LIFT improves sample efficiency in regression tasks.** The experiments are conducted on Insurance and Student datasets. The second column indicates the fraction of samples used for training the model. We observe no significant improvements in the performance when feature names are properly included.

Dataset (ID)	Corruption	ODC	LogReg	KNN	DT	MLP	SVM	RF	XG	LIFT/GPT-3	LIFT/GPT-J
Blobs (2)	0%	25.00	96.75	95.50	97.00	97.00	96.75	97.00	96.00	96.58	96.17
	5%	25.00	97.00	95.25	95.75	97.25	96.50	96.75	95.50	96.08	94.83
	10%	25.00	96.50	94.50	95.00	97.00	96.75	96.75	95.25	95.08	91.38
	20%	25.00	93.75	90.00	95.50	97.50	97.00	96.75	94.25	93.83	83.12
Moons (4)	0%	50.00	88.75	100.00	99.25	99.75	100.00	99.75	99.75	99.83	99.58
	5%	50.00	89.25	100.00	97.25	100.00	100.00	99.75	99.50	99.08	96.50
	10%	50.00	89.50	99.00	99.00	99.50	100.00	99.50	99.25	98.50	94.00
	20%	50.00	89.50	94.75	94.75	99.50	100.00	99.50	92.00	94.25	79.88
CMC (23)	0%	42.71	50.51	52.20	56.61	56.27	57.63	52.88	52.54	57.18	49.83
	5%	42.71	49.15	50.17	57.29	57.29	55.93	53.90	54.92	55.82	50.28
	10%	42.71	49.49	46.44	55.59	55.93	54.24	53.90	56.95	57.06	48.47
	20%	42.71	47.12	46.78	52.54	55.25	54.92	50.17	48.14	55.71	45.42
TAE (48)	0%	35.48	51.61	61.29	67.74	58.06	61.29	77.42	64.52	50.54	61.29
	5%	35.48	54.84	61.29	67.74	45.16	67.74	64.52	74.19	45.16	53.76
	10%	35.48	41.94	45.16	54.84	32.26	32.26	51.61	54.84	52.69	46.24
	20%	35.48	29.03	48.39	48.39	32.26	45.16	48.39	45.16	47.31	35.48
Pollen (871)	0%	50.00	49.09	46.88	48.96	49.22	51.56	45.97	48.31	49.57	50.39
	5%	50.00	51.43	48.18	49.22	50.26	49.74	48.44	46.62	50.65	48.61
	10%	50.00	48.70	48.70	47.27	50.00	46.49	51.17	47.01	48.96	48.66
	20%	50.00	50.39	49.22	50.52	47.01	50.52	47.14	49.61	50.74	50.82
Climate (1467)	0%	91.67	89.81	89.81	91.67	91.67	87.96	91.67	90.74	91.67	87.04
	5%	91.67	89.81	91.67	87.04	91.67	91.67	90.74	87.96	91.36	85.49
	10%	91.67	90.74	88.89	90.74	88.89	91.67	91.67	88.89	91.67	83.80
	20%	91.67	90.74	81.48	83.33	88.89	91.67	89.81	87.04	91.67	76.39
LED (40496)	0%	12.00	70.00	69.00	64.00	72.00	67.00	63.00	65.00	68.00	67.33
	5%	12.00	72.00	69.00	68.00	72.00	65.00	73.00	67.00	68.33	60.33
	10%	12.00	72.00	61.00	66.00	73.00	67.00	69.00	65.00	70.33	56.00
	20%	12.00	70.00	43.00	63.00	69.00	64.00	69.00	62.00	66.67	47.00

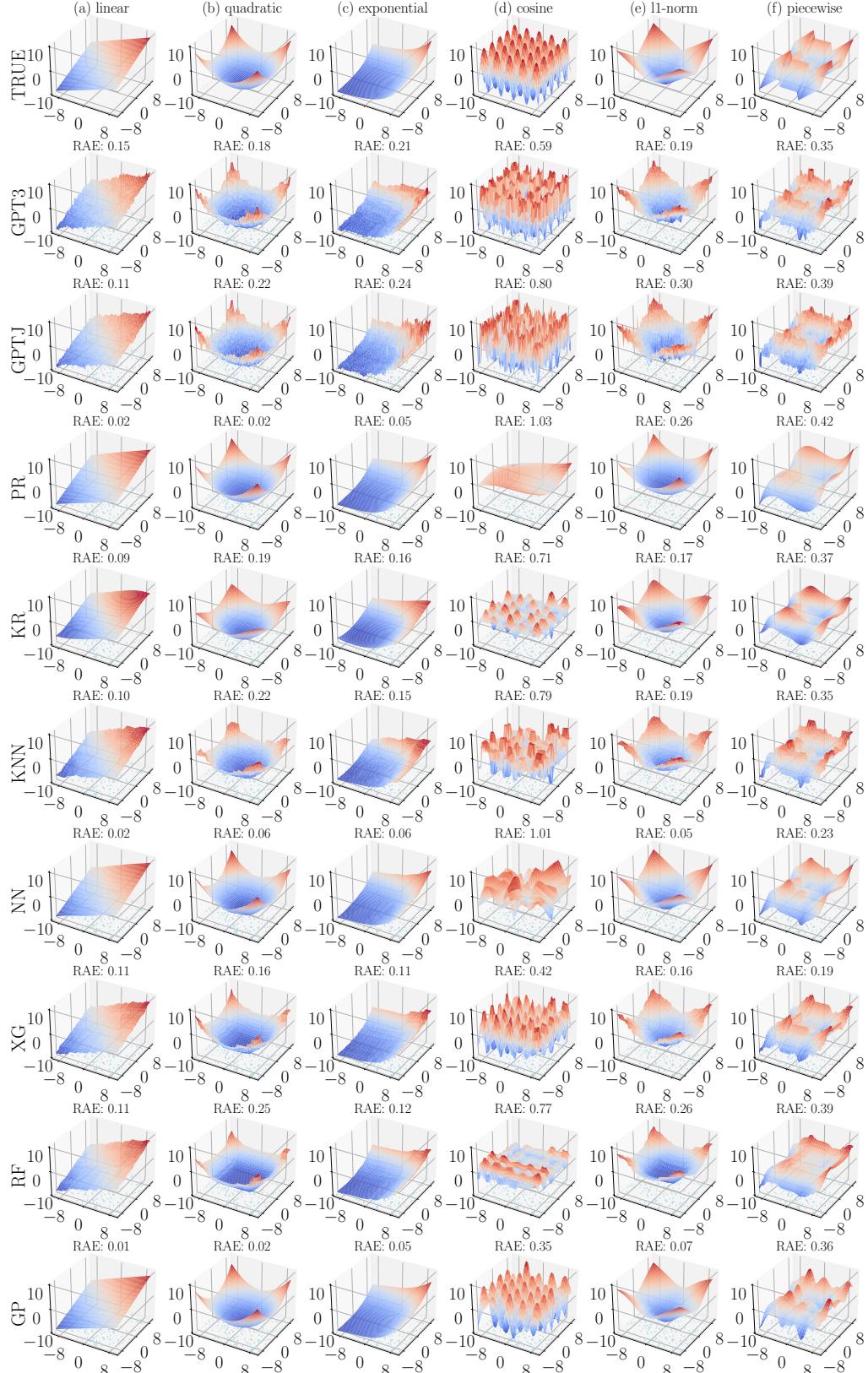
**Table 34: Accuracy ( $\uparrow$ ) comparison of various methods fitted to randomly corrupted classification labels.** In this regime, we corrupt a sample by assigning it another random label in the label space.

### B.2.6 Visualization of Regression Models

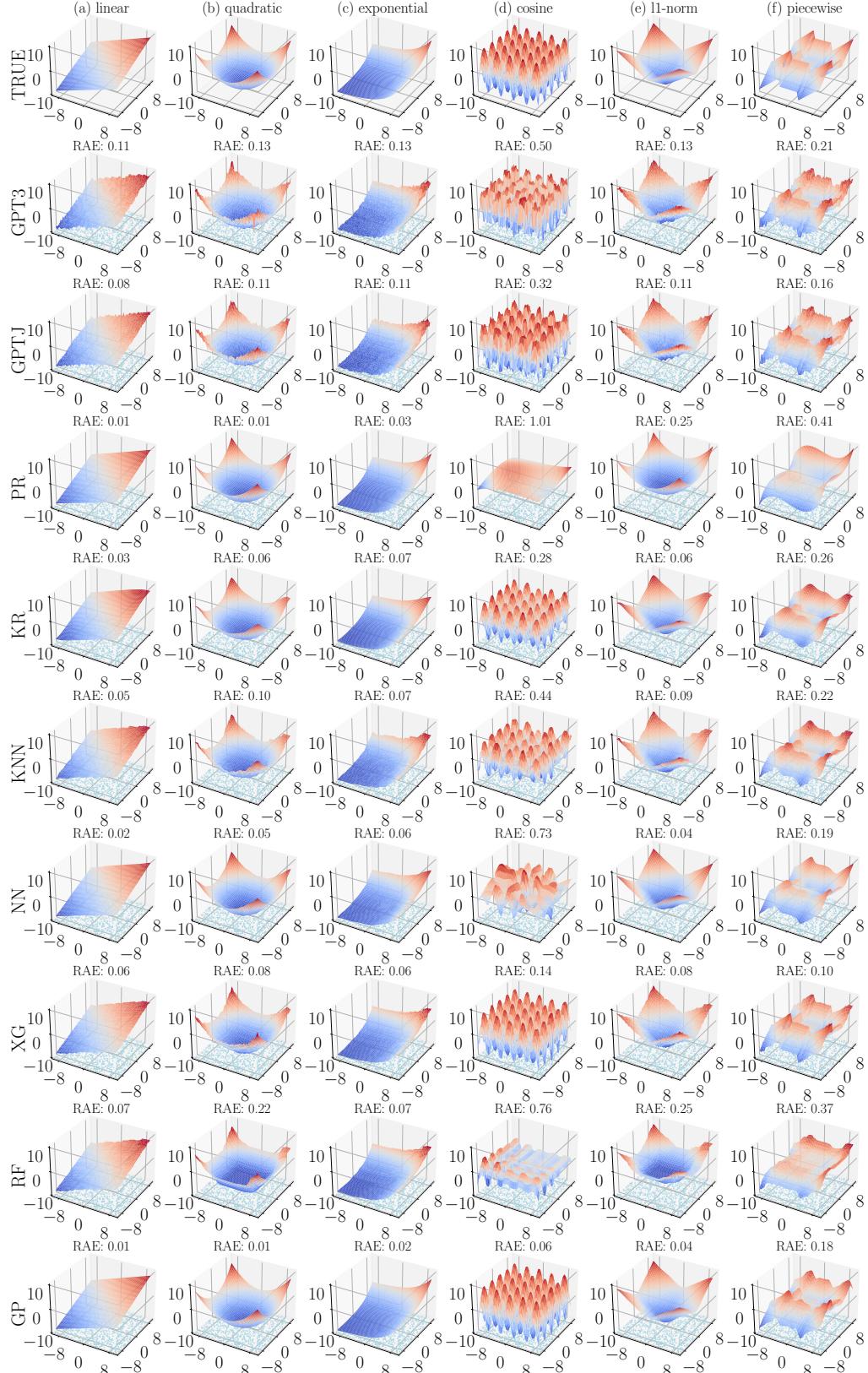
Fig. 36 and 37 visualize the 2D predictions for various functions with 200 and 1000 samples training datasets, respectively. Each coordinate of the training sample is drawn uniformly from  $[-10, 10]$ . Specifically, the prediction is performed on the interval  $[-12, 12]$ .

### B.3 Can LIFT Perform Ridge Regression via Data Augmentation?

As shown in Fig. 1, LIFT can perform linear regression. We take one step further and study whether LIFT can perform Ridge regression. Note that this is a non-trivial task as the LIFT framework does not allow any changes to the loss function. Consider a standard ridge regression problem solving the optimal  $\mathbf{w}$  with  $p$  parameters so that  $\|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \lambda\|\mathbf{w}\|_2^2$  is minimized. Note that this problem is equivalent to minimizing  $\|[\mathbf{y}^T, 0]^T - [\mathbf{X}^T, \sqrt{\lambda}\mathbf{I}]^T\mathbf{w}\|_2^2$ . Therefore, if we add  $p$  additional training samples  $\sqrt{\lambda}\mathbf{I}$ , one can perform ridge regression via data augmentation. Inspired by this, we study whether one can perform ridge regression via data augmentation within the framework of the LIFT framework. The results of LIFT on Ridge regression are reported in Table 38.



**Figure 36: Performance of LIFT/GPTs and baselines in approximating various functions.** The first row visualizes the true values of the functions, and the second & third rows visualize the predicted values of LIFT/GPTs after finetuning for the corresponding regression tasks with **200 training samples**. We compared with other baselines with the same training samples.



**Figure 37: Performance of LIFT/GPTs in approximating various functions.** The first row visualizes the true values of the functions, and the second & third rows visualize the predicted values of LIFT/GPTs after finetuning for the corresponding regression tasks with **1000 training samples**. We compared with other baselines with the same training samples.

Dataset ID	Corruption	ODC	LogReg	KNN	DT	MLP	SVM	RF	XG	LIFT/GPT-3	LIFT/GPT-J
Blobs (2)	0%	25.00	96.75	95.50	97.00	97.00	96.75	97.00	96.00	96.50	96.17
	5%	25.00	94.25	95.50	94.75	97.50	96.75	96.50	94.25	96.25	94.75
	10%	25.00	90.75	94.25	94.75	96.75	96.75	96.00	95.00	94.92	90.17
	20%	25.00	85.25	87.00	97.00	94.50	96.75	96.50	91.50	92.58	81.07
LED (40496)	0%	12.00	70.00	69.00	64.00	72.00	67.00	63.00	65.00	65.67	67.33
	5%	12.00	71.00	69.00	67.00	67.00	68.00	64.00	65.00	69.33	58.00
	10%	12.00	72.00	64.00	70.00	70.00	64.00	68.00	65.00	70.00	55.67
	20%	12.00	70.00	63.00	67.00	74.00	66.00	66.00	65.00	63.33	53.00

**Table 35: Accuracy ( $\uparrow$ ) of various methods fitted to systematically corrupted classification labels.** In this regime, we corrupt a label by assigning all corrupted labels of one class to a single label.

$p$	$\lambda$	LIFT/GPT-J		LIFT/GPT-3	
		LR	Ridge	LR	Ridge
1	0	$0.000 \pm 0.000$	$0.000 \pm 0.000$	$0.915 \pm 0.000$	$0.000 \pm 0.000$
1	10	$0.000 \pm 0.000$	$0.016 \pm 0.000$	$0.915 \pm 0.000$	$0.016 \pm 0.000$
1	50	$0.000 \pm 0.000$	$0.403 \pm 0.000$	$0.915 \pm 0.000$	$0.402 \pm 0.000$
1	100	$0.000 \pm 0.000$	$1.691 \pm 0.000$	$0.915 \pm 0.000$	$1.690 \pm 0.000$
1	1000	$0.000 \pm 0.000$	$170.406 \pm 0.000$	$0.915 \pm 0.000$	$170.612 \pm 0.000$
10	0	$0.532 \pm 0.000$	$0.532 \pm 0.000$	$0.915 \pm 0.000$	$0.521 \pm 0.000$
10	10	$0.374 \pm 0.000$	$0.369 \pm 0.000$	$0.915 \pm 0.000$	$0.504 \pm 0.000$
10	50	$0.417 \pm 0.000$	$0.523 \pm 0.000$	$0.915 \pm 0.000$	$0.563 \pm 0.000$
10	100	$0.365 \pm 0.000$	$1.307 \pm 0.000$	$0.915 \pm 0.000$	$1.539 \pm 0.000$
10	1000	$0.414 \pm 0.000$	$114.042 \pm 0.000$	$0.915 \pm 0.000$	$111.357 \pm 0.000$
50	0	$0.688 \pm 0.000$	$0.688 \pm 0.000$	$0.915 \pm 0.000$	$1.064 \pm 0.000$
50	10	$0.628 \pm 0.000$	$0.635 \pm 0.000$	$0.915 \pm 0.000$	$0.909 \pm 0.000$
50	50	$0.553 \pm 0.000$	$0.732 \pm 0.000$	$0.915 \pm 0.000$	$1.296 \pm 0.000$
50	100	$0.774 \pm 0.000$	$1.857 \pm 0.000$	$0.915 \pm 0.000$	$2.311 \pm 0.000$
50	1000	$0.970 \pm 0.000$	$118.241 \pm 0.000$	$0.915 \pm 0.000$	$133.122 \pm 0.000$

**Table 38: Performance of LIFT on Ridge regression.** We measure the RAE ( $\downarrow$ ) of LIFT corresponding to Linear Regression (LR) and Ridge Regression. The RAEs indicate that LIFT does not perform well on the Ridge regression problem.

We observe that LIFT fails to perform Ridge regression. This is expected, as LIFT is shown to be robust to outliers (in Sec. 3.4).