
DEALING OUT OF DISTRIBUTION WITH PREDICTION PROBLEM

Anonymous authors

Paper under double-blind review

ABSTRACT

The open world assumption in model development means that a model may lack sufficient information to effectively handle data that is completely different or out of distribution (OOD). When a model encounters OOD data, its performance can significantly decrease. Improving the model’s performance in dealing with OOD can be achieved through generalization by adding noise, which can be easily done with deep learning. However, many advanced machine learning models are resource-intensive and designed to work best with specialized hardware (GPU), which may not always be available for common users with hardware limitations. To provide a deep understanding and solution on OOD for general user, this study explores detection, evaluation, and prediction tasks within the context of OOD on tabular datasets using common consumer hardware (CPU). It demonstrates how users can identify OOD data from available datasets and provide guidance on evaluating the OOD selection through simple experiments and visualizations. Furthermore, the study introduces Tabular Contrast Learning (TCL), a technique specifically designed for tabular prediction tasks. While achieving better results compared to heavier models, TCL is more efficient even when trained without specialised hardware, making it useful for general machine-learning users with computational limitations. This study includes a comprehensive comparison with existing approaches within their best hardware setting (GPU) compared with TCL on common hardware (CPU), focusing on both accuracy and efficiency. The results show that TCL exceeds other models, including gradient boosting decision trees, contrastive learning, and other deep learning models, on the classification task.

1 INTRODUCTION

The concept of open-world assumption in model development means that a model may not have enough information to effectively handle data that is completely different or out of distribution (OOD). When a model meets OOD data, it may suffer a significant decrease in performance (Hsu et al., 2020; Hendrycks and Gimpel, 2016). To handle this, model generalisation by introducing noise can be used, which can be achieved easily with deep learning. However, advanced deep learning algorithms such as FT-Transformer benefit from the advancement of specialised hardware such as GPU or TPU (Hwang, 2018), this type of hardware is not always available to the general user (Ahmed and Wahed, 2020). These demand the user to find the best way to deal with these challenges and emphasize the importance of our research.

While out-of-distribution (OOD) detection has been extensively studied (Lee et al., 2020a), the challenge of prediction tasks for OOD data, particularly in tabular datasets, remains underexplored. Significant progress has been made in OOD detection with algorithms like MCCD (Lee et al., 2020b), OpenMax (Bendale and Boult, 2016), Monte Carlo Dropout (Gal and Ghahramani, 2016), and ODIN (Liang et al., 2017). However, the study of prediction tasks on OOD for tabular data is limited. Tree-based classical models are known to

047
048
049
050
051
052
053
054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093

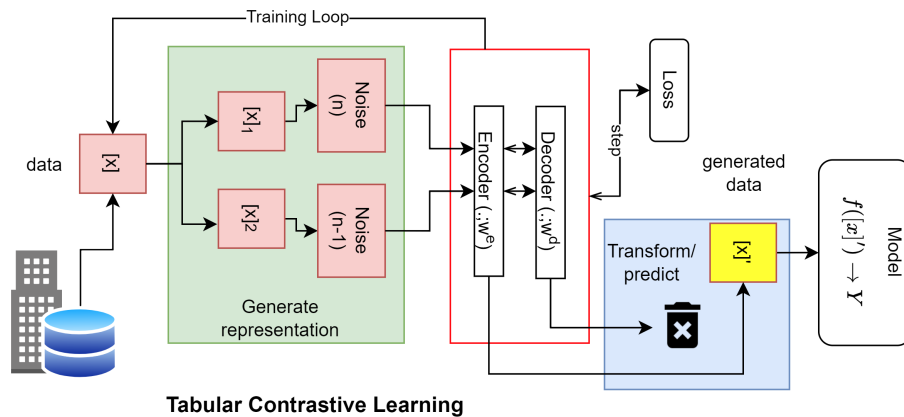


Figure 1: Tabular Contrastive Learning (TCL). The data $[x]$ is duplicated ($[x]_1, [x]_2$) and noise is added. Both duplicates are then encoded and decoded to compute the loss. During inference, TCL only uses the encoder to produce new data $[x']$ that enhances supervised learning performance $f([x']) \rightarrow Y$. The decoder is omitted during inference and used only for training.

be reliable for tabular data (Grinsztajn et al., 2022), but our experiments show that these models exhibit a decrease in performance when dealing with OOD data.

In this study, we made several contributions. First, we show step by step how to implement existing methods for detecting, separating, evaluating, and visualizing out-of-distribution (OOD) data using real-world datasets. Second, we assess the performance of existing tabular machine learning algorithms in handling OOD data. Lastly, we introduce a new approach called TCL, which provides efficiency and flexibility while achieving comparable performance.

Tabular Contrast Learning (TCL), Figure 1, is a local adaptation of Contrastive Federated Learning (CFL) (Ginanjari et al., 2024) designed for prediction tasks on tabular datasets for a general user. TCL is based on the principles of contrastive learning (Chen et al., 2020; Ucar et al., 2021) but is optimized for tabular data structures. TCL approach offers several advantages, e.g. **Efficiency**: TCL is designed to be faster and more compact compared to current state-of-the-art models, **Flexibility**: TCL can be integrated with various supervised learning algorithms and **Performance**: TCL achieves competitive performance.

Our experiment demonstrates that TCL delivers performance and efficiency (Huang et al., 2017) (defined by a higher speed/accuracy trade-off score) compared to other models.

2 RELATED WORK

2.1 OOD DETECTION

OpenMax (Bendale and Boult, 2016) uses the concept of Meta-Recognition to estimate the probability that an input belongs to an unknown class. OpenMax characterizes the failure of the recognition system and handles unknown/unseen classes during operation. In deep learning, SoftMax calculates as $P(y = j|x) =$

$$\frac{e^{v_j(x)}}{\sum_{i=1}^N e^{v_i(x)}}$$

OpenMax recognizes that in out of distribution (OOD), the denominator of the SoftMax layer does not require the probabilities to sum to 1.

094 **TemperatureScaling** (Platt et al., 1999) is a single-parameter variant of the Platt scaling. In a study by Guo
095 *et al.* (Guo et al., 2017), despite its simplicity, temperature scaling is effective in calibrating a model for
096 deep learning. This also suggests that temperature scaling can be used to detect OOD. Our study uses these
097 two approaches to separate OOD from the dataset and use it as validation data.

098 **Multi-class classification, deep neural networks, Gaussian discriminant analysis (MCCD)** (Lee et al.,
099 2020b) is OOD detection algorithm based deep neural network that claim to have better classification infer-
100 ence performance. It is focuses on finding sperical-decission across classes.

101 Our work mainly uses OpenMax and TemperatureScaling. While the original algorithms are not new, both
102 algorithm have latest update and better support under pytorch-ood (Kirchheim et al., 2022) compared to
103 MCCD (Lee et al., 2020c).

106 2.2 TABULAR DATA PREDICTION

107
108 **Neural Network-based Methods:** Multilayer Perceptron (MLP) (Ruck et al., 1990; Gorishniy et al., 2023):
109 A straightforward deep learning approach for tabular data. Self-Normalizing Neural Networks (SNN).
110 (Klambauer et al., 2017): Uses SELU activation to train deeper networks more effectively.

111 **Advanced Architectures:** Feature Tokenizer Transformer / FT-Transformer (Vaswani et al., 2017): Adapts
112 the transformer architecture for tabular data, consistently achieving high performance. Residual Network
113 / ResNet (Li et al., 2018): Utilizes parallel hidden layers to capture complex feature interactions. Deep &
114 Cross Network / DCN V2 (Wang et al., 2020): Incorporates a feature-crossing module with linear layers and
115 multiplications. Automatic Feature Interaction / AutoInt (Song et al., 2018): Employs attention mechanisms
116 on feature embeddings. Neural Oblivious Decision Ensembles / NODE (Popov et al., 2019): A differentiable
117 ensemble of oblivious decision trees. Tabular Network / TabNet (Arık and Pfister, 2019): Uses a recurrent
118 architecture with periodic feature weight adjustments. Focuses on attention framework.

119 **Ensemble Methods:** GrowNet (Badirli et al., 2020): Applies gradient boosting to less robust MLPs, pri-
120 marily for classification and regression tasks.

121 **Gradient Boosting Decision Tree (GBDT)** (Grinsztajn et al., 2022) : XGBoost :A tree-based ensemble
122 method that uses second-order gradients and regularization to prevent overfitting while maximizing compu-
123 tational efficiency. LightGBM :A fast and memory-efficient boosting framework that uses histogram-based
124 algorithms and leaf-wise tree growth strategy for faster training. CatBoost :A gradient boosting implementa-
125 tion specifically optimized for categorical features with built-in ordered boosting to reduce prediction shift.

126 These models have shown varying degrees of success in tabular data prediction. However, their performance
127 on OOD data remains a critical area for investigation. We include mentioned model as our base models.

130 2.3 TABULAR CONTRASTIVE LEARNING

131 **SubTab** (Ucar et al., 2021) and **SCARF** (Bahri et al., 2022) are a contrastive learning model tailored for
132 tabular datasets. Similar to the fundamental concept of contrastive learning for the image of SimCLR (Chen
133 et al., 2020). SubTab and SCARF calculate contrastive loss using cosine or euclidean distance. **CFL** (Gi-
134 nanjar et al., 2024) is a federated learning algorithm proposed to tackle vertical partition within data silos.
135 CFL explores the possibility of implementing contrastive learning within vertically partitioned data without
136 the need for data sharing. CFL merge the weight by understanding that the data came from global imaginary
137 dataset which is vertically partitioned. CFL uses contrastive learning as a medium for black box learning.
138 CFL focuses on collaborative learning across silos. In this study, we study learning from local data with
139 OOD, a problem that is yet to be explored by CFL CFL focuses on a Federated learning network, while ours
140 is common tabular data. CFL, while exhibiting a similar name, uses partial data augmentation as part of

141 the federated learning concept and is similar to image contrastive learning. TCL, in the other hand use full
142 matrix augmentation to support tabular data.

143 144 145 3 PROBLEM FORMULATION

146 147 148 3.1 DEFINITION

149 **Definition 1: Tabular Data.** Let $D \in R^{nd}$ be a tabular dataset with n samples and d features, and $Y \in R^n$
150 be the corresponding labels. Tabular data is characterized by its structured format, where each row represents
151 a sample and each column represents a feature.

152 **Definition 2: Out-of-Distribution Prediction.** Given a model trained on in-distribution data $D_{in} =$
153 $\{(x_i, y_i) | x_i \in X_{in}, y_i \in Y_{in}\}$, where X_{in} follows a distribution p_{in} , the task is to make accurate predictions
154 on OOD data X_{ood} that follows a different distribution p_{ood} , where $p_{ood} \neq p_{in}$.

155 **Definition 3: Efficiency-Accuracy Trade-off.** We use a speed/accuracy trade off (Huang et al., 2017) from
156 the total performance matrix. Let T be the total performance, then $T = \frac{P}{t}$ for classification or $T = \frac{1/P}{t}$ for
157 regression.

158
159 A performance evaluation P used is F1 or RMSE for the prediction task and the time t in seconds for the
160 duration of training. The P is used to obtain the standard performance of a model. The t is used to evaluate
161 the time it takes to train a model. A smaller t means a smaller resource to find the best model by tuning the
162 hyperparameters. The adjustment $\frac{1}{P}$ for regression is necessary because in regression tasks, RMSE is used
163 as the performance metric, and smaller values are better.

164 165 166 3.2 PROBLEM STATEMENT

167
168 In machine learning, the goal is to build a model $f : x \rightarrow y$ that generalizes unseen data. However, if the
169 model is exposed to samples outside the distribution during inference, it may make unreliable predictions or
170 exhibit unexpected behaviour.

171 Formally, a prediction task can be defined as finding a model $f : (\cdot)$ that minimizes the expected error over
172 a dataset D_{in} with distribution p_{in} . This can be defined as $\mathbf{min} \text{Error}(x, y)_D = [L(f(x), y)]$, where L is a
173 loss function that measures the discrepancy between the prediction of the model $f(x)$ and the true label y . A
174 bigger E means poor model performance P or can be denoted as $E \uparrow = P \downarrow$. When a different distribution
175 p_{ood} is introduced to a model, the performance decreased or $P(f((x)_{D_{in}})) > P(f((x)_{D_{ood}}))$.

176 The time t to find the best model should also be considered. A time-consuming model takes more resources
177 to train and tune. When dealing with a large dataset, the time used to train and tune a model is a big concern.
178 A larger and longer model does not always equal better performance P . This can be written as $t \uparrow \neq P \uparrow$.

179 We evaluated the total/overall performance of the models when dealing with tabular dataset with OOD. The
180 overall objection can be written as $T = \frac{\max P_{ood}}{\min t} = \frac{1/\min E(x,y)_{D_{ood}}}{\min t}$.

181 182 183 4 PROPOSED METHOD

184
185
186 We introduce Tabular Contrastive Learning (TCL), an improved approach designed to enhance prediction
187 tasks on tabular data, particularly with hardware limitations.

4.1 TCL MAIN ARCHITECTURE

Augmented Data

As a contrastive learning based algorithm, TCL works based on augmented data. TCL creates two data augmentations. Denoted as $\{x^1, x^2\} = \text{Aug}(x)$. Noise was added to these augmented data.

Matrix Augmentation

In TCL, all original data (without slicing or splitting) is utilized as a representation. Unlike previous approaches such as simCLR (Chen et al., 2020) and SubTab (Ucar et al., 2021) that use slice, TCL utilizes the entire original data matrix for representation. This allows for capturing more comprehensive feature interactions in tabular data.

Encoder-Decoder Structure

The contrastive learning architecture includes two main components: an encoder and a decoder. The encoder transforms input data into a compressed representation called the latent space, denoted $E : x; \omega^e \rightarrow x^e$ where ω is parameter and e is encoder notation, while the decoder reconstructs the original data from this representation denoted $P : x^e; \omega^p \rightarrow x^p$ where p is notation for decoder. During training, both components are used, but only the encoder is employed during inference, allowing efficient compression of new data into its learned latent representation without reconstruction.

Modified Contrastive Loss

Loss is calculated based on augmented data, not original data. TCL simplifies contrastive loss calculation to enhance both performance and training speed. In contrastive learning, the loss is computed based on the similarity or dissimilarity of the augmented noisy data. Since TCL deals with row-based tabular data and it is a unsupervised learning, the method used is similarity. TCL aims to pull the noisy data points that originated from the same data. This is achieved by minimizing the total loss L_c between representations. During training, the total loss is calculated as follows:

$$L_t(x) = (L_r(x) + L_c(x) + L_d(x)) \quad (1)$$

Where L_r is the reconstruction loss, L_c is the contrastive loss, and L_d is the distance loss. The objective of contrastive learning is to minimize the total loss L_t . When D is dataset, and sliced data $\mathcal{B} \in D$, then:

$$\min L_t(.; \omega^e, \omega^p) = \min \frac{1}{J} \sum_{j=1}^J L_t(P(E(.; \omega^e); \omega^p)) \quad (2)$$

with w is weight, $P(\cdot)$ is decoder function, and $E(\cdot)$ is encoder function. When $MSE(\cdot)$ is the mean square error function, and $[x]$ is noisy data of \mathcal{B} , then: $L_r(x) = \frac{1}{N} \sum_n^N MSE(\hat{x}, x)$ and $L_d(x) = \frac{1}{N} \sum_n^N MSE(x^{e1}, x^{e2})$. While L_r is calculated from decoded data and original data, L_d is calculated from encoded data only. We simplified the contrastive loss L_c by using only the result of a dot product compared with other contrastive learning that used euclidean distance.

$$L_c(x) = \frac{1}{N} \sum_n^N \left(-\log \frac{\exp(MSE([0], \text{dot}(x^{e1} \cdot x^{e2}))/T))}{\sum_{k=1}^K \exp(MSE([0], \text{dot}(x^{e1} \cdot x^{e2}))/T)} \right) \quad (3)$$

5 TCL ALGORITHM

The TCL process involves several steps. First, a minibatch of N samples is sampled from the dataset. Then, for each sample in the batch, two augmented views are created and added with noise. Although they came

235 from the same data, both augmented data are different due to previous treatments. These augmented views
236 are passed through an encoder network to obtain encoded representations. A decoder is then applied to
237 the encoded representations. The loss function then calculates the difference between two augmented data.
238 By minimizing this loss, noisy data is pulled together. Because TCL applies full matrix representation, the
239 process pulls noisy data row by row together. This results in generalized data for better inference. The
240 complete steps are presented on Algorithm 2 in the Appendix section.

241 242 243 6 EXPERIMENT

244 245 246 6.1 DATASETS

247
248 We utilized 10 diverse tabular datasets. The datasets are Adult (Becker and Kohavi, 1996), Helena (Guyon
249 et al., 2019), Jannis (Guyon et al., 2019), Higgs Small (Baldi et al., 2014), AloI (Geusebroek et al., 2005),
250 Epsilon (PASCAL Challenge on Large Scale Learning, 2008), Cover Type (Blackard and Dean, 2000),
251 California Housing (Pace and Barry, 1997), Year (Bertin-Mahieux et al., 2011), Yahoo (Chapelle and Chang,
252 2011), and Microsoft (Qin and Liu, 2013).

253 254 255 6.2 OOD DETECTION

256
257 We have implemented two Out-of-Distribution (OOD) detection methods, namely OpenMax (Bendale and
258 Boulton, 2016) and TemperatureScaling (Platt et al., 1999). We applied OOD detection methods to each
259 dataset to transform the data and establish thresholds. The thresholds were manually assigned by observing
260 the graphs produced with the OOD detection algorithm. The manual assignment was done by selecting a
261 single point on a tail of the observation. We then separated the OOD data based on the thresholds to generate
262 two sets, D_{in}^M and D_{ood}^N , where M and N are total sample in each set, ($D_{in} + D_{ood} = D$). We expect to
263 find $M > N$. The OOD separation is validated using linear regression. Finally, we compared the model
264 performance with and without OOD separation, expecting to find that the performance decreased in this step
265 $f : D_{in} > f : D_{ood}$.

266 The results of the experiment is two set of dataset D_{in}, D_{ood} . While D_{in} is used for train dataset, D_{ood} is
267 used for test dataset. The Algorithm 1 in the Appendix section explains step by step the OOD detection is
268 done.

269 270 271 272 6.3 PREDICTION ON OOD DATASET

273
274 We experimented with 12 based models. For deep learning tabular models we use FT-T, DCN2, GrowNet,
275 ResNet, MLP, AutoInt and TabR-MLP (Gorishniy et al., 2023). In addition, we experimented with the recent
276 implementation of contrastive learning for tabular data SubTab (Ucar et al., 2021) and SCARF (Bahri et al.,
277 2022), which comes from a similar domain to our TCL. We also did not apply FT-Transformer to some
278 datasets. The FT-Transformer is heavy and has reached our hardware limitation. Finally, we compared our
279 TCL with GDBT models.

280 Our experiment used an NVIDIA H100 GPU for all models except TCL and the GDBT-based model. TCL
281 and GDBT were trained on a CPU (Apple / AMD) to emphasize our advancement within limited hardware.

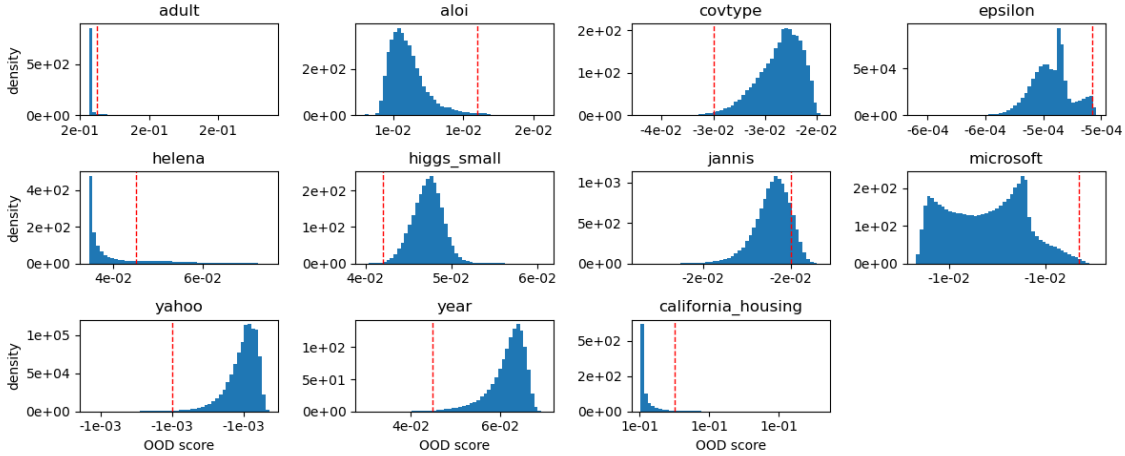


Figure 2: Histogram representing the OOD scores across various datasets. The red line is the threshold line that indicates whether the data is Out of Distribution or not.

7 RESULT AND EVALUATION

7.1 OOD DETECTION

Table 1: The OOD detection settings. Performances are results of model trained with linear regression (r^2) and logistic regression (accuracy). When OOD dataset is separated and used as test dataset in (^b) the performance of the model is decreased. OOD case in the epsilon (^{*}) dataset cannot be identified.

Dataset	Detectors	Norms	OOD threshold	Accuracy without OOD ^a		Accuracy with OOD ^b	
				Train	Test	Train non-OOD	Test OOD
Adult	OpenMax	12	0.162800	0.782561	0.783089	0.797838	0.267408
Helena	OpenMax	11	0.045000	0.194532	0.196626	0.146428	0.047553
Jannis	TemperatureScaling	11	-0.020000	0.561953	0.563982	0.577274	0.474213
Higgs small	OpenMax	11	0.042000	0.622216	0.616879	0.622667	0.568493
Aloi	OpenMax	11	0.016000	0.260648	0.232546	0.328289	0.071186
Epsilon [*]	TemperatureScaling	12	-0.000523	0.898635	0.897230	0.898700	0.893338
Covtype	TemperatureScaling	11	-0.035000	0.604674	0.604373	0.603514	0.435439
California Housing	OpenMax	11	0.110000	0.333448	0.180680	0.477136	-6.033595
Year	OpenMax	12	0.045000	0.168659	0.167098	0.169745	-0.714197
Yahoo	TemperatureScaling	11	-0.001460	0.325685	0.326275	0.325577	-0.297845
Microsoft	TemperatureScaling	11	-0.008300	0.045648	0.044109	0.047334	-0.841856

Table 1 shows significant differences between the two settings. Without OOD (Table 1, Section a), the training and test results are comparable. However, when used as test data, the OOD reduces the performance of the models. OOD leads to a 20% decrease (Table 1, Section b) in performance between training and test results for the classification task, and a negative r^2 for the regression task. The OOD separation process is

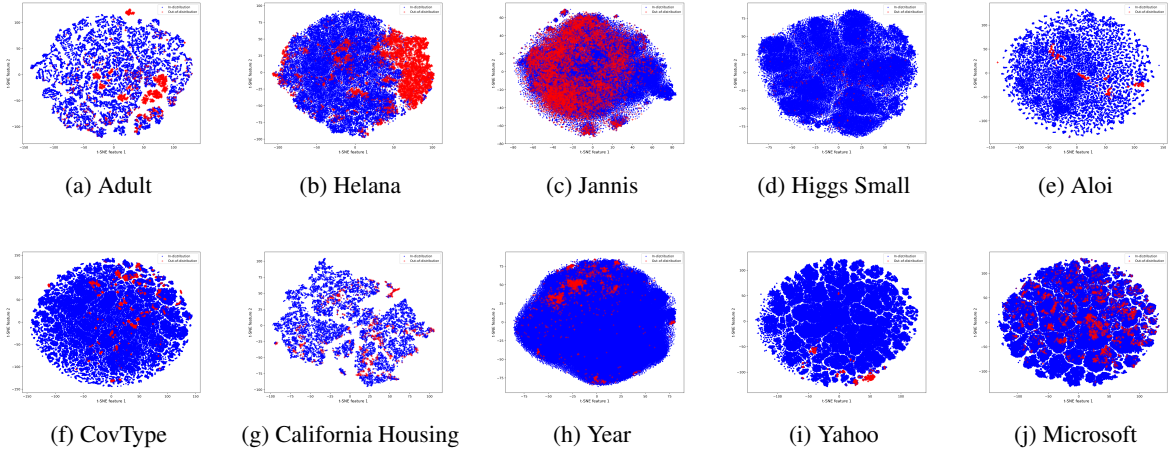


Figure 3: OOD Visualisation. (.) or blue color is the in distribution data (ID), (+) or red is out of distribution data (OOD). In all dataset, except for Jannis, it is clear that OOD fills empty space between ID. Jannis data is evenly distributed, the visualisation capture neither ID nor OOD.

Table 2: Experiment result. F1 score for classification and RMSE for regression. Datasets with (*) mean a regression problem. Models with (c) are contrastive learning based models.

	AD \uparrow	HE \uparrow	JA \uparrow	HI \uparrow	AL \uparrow	CO \uparrow	CA* \downarrow	YE* \downarrow	YA* \downarrow	MI* \downarrow
FT-T	0.782	0.153	0.572	0.738	0.407	-	0.867	6.461	-	-
DCN2	0.744	0.129	0.542	0.710	0.414	0.58	2.602	7.054	0.645	0.746
GrowNet	0.465	-	-	0.685	-	-	0.969	7.605	1.01	0.769
ResNet	0.652	0.10	0.574	0.753	0.437	0.694	0.892	6.496	0.639	0.736
MLP	0.508	0.146	0.561	0.753	0.326	0.617	0.894	6.488	0.657	0.741
AutoInt	0.78	0.133	0.549	0.719	0.401	0.608	0.89	6.673	-	0.739
TabR-MLP	0.688	0.165	0.541	0.753	0.429	0.688	2.677	2e5	1.285	0.79
TCL ^c	0.831	0.154	0.575	0.758	0.447	0.880	0.843	6.491	0.652	0.738
Scarf ^c	0.720	0.00	0.122	0.308	0.00	0.091	-	-	-	-
SubTab ^c	0.714	0.146	0.504	0.602	0.322	0.59	1.012	6.668	0.656	0.744

visualized in Figure 2, which presents histogram graphs of the transformed data with detectors. The Epsilon dataset shows two peaks, indicating complexity, and outliers cannot be detected. In contrast, the Microsoft dataset shows a performance decrease with OOD. Figures 3 show the results of separating in-distribution (ID) and out-of-distribution (OOD) data using TSNE in an unsupervised manner. The figures demonstrate that OOD data fills the empty space between ID data, indicating a strong presence of OOD. Table 1, Figure 2, and Figure 3 show strong indications of the existence of out-of-distribution (OOD) data.

7.2 MODELS PERFORMANCE

Table 2 shows the results of the experiment. Overall, TCL outperforms other models, while the performance of the other models is comparable across various datasets. There are some exceptions where specific models underperform relative to others. For instance, GrowNet performs below average on the adult dataset, DCN V2 underperforms on the California Housing dataset, and GrowNet also underdelivers on the Yahoo dataset.

Table 3: Experiment result of TCL compared to GBDT. F1 score for classification and RMSE for regression. Datasets with (*) mean a regression problem. Model with (c) means contrastive learning based model, models with (x) mean GBDT based models.

	AD↑	HE↑	JA↑	HI↑	AL↑	CO↑	CA*↓	YE*↓	YA*↓	MI*↓
TCL ^c	0.831	0.154	0.575	0.758	0.447	0.880	0.843	6.491	0.652	0.738
Lightgbm ^x	0.591	0.080	0.432	0.609	0.177	0.219	0.848	6.565	0.661	0.740
CatBoost ^x	0.927	0.152	0.533	0.718	-	0.753	0.827	6.622	0.655	0.733
XGB ^x	0.925	0.127	0.532	0.739	0.328	0.700	0.845	6.867	0.654	0.739

Table 4: Table of training duration in second of each dataset. Datasets with (*) means a regression problem. All model except for TCL are trained with GPU. TCL were trained with CPU

	AD↓	HE↓	JA↓	HI↓	AL↓	CO↓	CA*↓	YE*↓	YA*↓	MI*↓
FT-T	1027	130	155	94	1205	-	88	1290	-	-
ResNet	2.1e+02	32	21	56	44	618	15	236	284	950
TCL	15	23	23	38	40	330	7	240	620	820

In contrast, TCL stands out by outperforming other models on most datasets, particularly in classification problems. Nevertheless, TCL’s performance in regression problems is not significantly behind that of the top models..

When compared with the GBDT method, TCL outperforms in most datasets, especially in the classification problem, see Table 3. Compared to other classification problems, the adult dataset has a relatively higher score across other datasets. This shows that adult dataset does not require generalisation during prediction, which also explains why TCL performs under CatBoost. CatBoost dominant in 4 datasets beat any other GBDT algorithm.

7.3 TRAINING DURATION

Table 4 displays the training duration for the best three deep learning models. Each model has unique characteristics and training steps, and all seven models (FT-T, DCN 2, GrowNet, ResNet, MLP, AutoInt, TabR-MLP) underwent extensive tuning. The Yahoo and Microsoft datasets required 5 days to complete the entire parameter-tuning process. For FT-T and resnet, a single training time was sampled once the tuning process was completed. TCL, which involve unsupervised training, The time recorded is time for each model to stabilize their loss with a 256 batch size, which is around 15 epochs. it is clear that TCL has a short training time.

7.4 EFFICIENCY EVALUATIONS

Table 5: A speed/accuracy trade off matrix $T = \frac{P}{t}$ where P performance matrix used and t is time in second required. A higher result is better. Datasets with (*) mean a regression problem.

	AD	HE	JA	HI	AL	CO	CA*	YE*	YA*	MI*
FT-T	0.00076	0.0012	0.0037	0.0079	0.00034	-	0.013	0.00012	-	-
ResNet	0.0031	0.0031	0.027	0.013	0.0099	0.0011	0.075	0.00065	0.0055	0.0014
TCL	0.055	0.0066	0.025	0.028	0.0199	0.0026	0.16	0.00064	0.0024	0.0016

Table 5 shows that TCLs are dominant. FT-Transformer and ResNet produce a good F1 and RMSE score; however, they take more time to train. In **FT-Transformer**, multiple attention heads process numeric and

423 categorical features separately before combining them. The model includes four types of layers that grow
 424 exponentially, leading to resource-intensive computations. **ResNet** employs parallel calculations across mul-
 425 tiple convolutional layers (SubNet), using three identical SubNets, one of which is highly filtered. In con-
 426 trast, TCL has a simpler architecture akin to MLP, achieving a high speed/accuracy trade-off. **TCL** features
 427 narrow layers for both the encoder and decoder, each with one hidden layer and one normalization layer, re-
 428 sulting in fewer layers than ResNet. However, TCL’s pair operation for loss calculation doubles its training
 429 time.

431 7.5 TCL EFFICIENCY EVALUATION

432 Our TCL has undergone significant algorithm modifications, making the original similarity loss function
 433 inapplicable. We compare TCL with SubTab, which employs a similarity function for tabular contrastive
 434 learning, to evaluate TCL’s efficiency.

435 Table 6 shows that the dot product applied on TCL consistently faster compared to similarity distance func-
 436 tion applied to similar contrastive learning under SubTab. The efficiency gain from using the dot product
 437 supports our decision to incorporate it into the TCL model. Implementation of the entire original data ma-
 438 trix for representation removes matrix splitting as used in common contrastive learning. Implementation of
 439 dot product removes the requirement to calculate more complex similarity scores. Both algorithms were
 440 evaluated under CPU.

441 Table 6: Table of training duration in second of each dataset. Datasets with (*) mean a regression problem.
 442 TCL uses the dot product, and SubTab uses the similarity function. Both algorithms were evaluated under
 443 CPU.

	AD↓	HE↓	JA↓	HI↓	AL↓	CO↓	CA*↓	YE*↓	YA*↓	MI*↓
TCL	15	23	23	38	40	330	7	240	620	820
SubTab	1400	1700	2400	2800	3400	1.8e+06	640	1.5e+04	3.6e+04	4e+04

450 8 CONCLUSION

451 The choice of models for tabular datasets with out-of-distribution (OOD) data depends on the user’s needs
 452 and available resources. TCL outperforms other heavier models for classification problems on OOD while
 453 maintaining efficiency. RestNet and FT-Transformer perform well on many datasets, but these models re-
 454 quire more resources, which may not always be feasible. It is worth noting that TCL was trained on a CPU,
 455 and RestNet and FT-T were trained on a GPU. This makes TCL available for more users than other models
 456 that require more training resources. Both RestNet and TCL can be options for fine-tuning and serving as
 457 head-to-head comparison models.

458 Although TCL has shown promising results, there are opportunities for potential enhancement. A continual
 459 learning can be proposed to improve performance. Further optimization of the contrastive learning process
 460 can be studied to achieve even greater efficiency. Additionally, there is a need to explore TCL’s perfor-
 461 mance on a wider range of domain-specific tabular datasets. Furthermore, it is crucial to investigate TCL’s
 462 interpretability, as this is important for many real-world applications.

465 9 REPRODUCIBILITY

466 The code for this work can be found online (anonymized, 2024) (submitted as a supplementary file). The
 467 dataset is also available online and can be downloaded using the information provided in the citation.

470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516

REFERENCES

- Yen-Chang Hsu, Yilin Shen, Hongxia Jin, and Zsolt Kira. Generalized odin: Detecting out-of-distribution image without learning from out-of-distribution data, 2020.
- Dan Hendrycks and Kevin Gimpel. A baseline for detecting misclassified and out-of-distribution examples in neural networks. *5th International Conference on Learning Representations, ICLR 2017 - Conference Track Proceedings*, 10 2016. URL <https://arxiv.org/abs/1610.02136v3>.
- Tim Hwang. Computational power and the social impact of artificial intelligence. *SSRN Electronic Journal*, 2018. doi: 10.2139/ssrn.3147971.
- Nur Ahmed and Muntasir Wahed. The de-democratization of ai: Deep learning and the compute divide in artificial intelligence research, 2020. URL <https://arxiv.org/abs/2010.15581>.
- Dongha Lee, Sehun Yu, and Hwanjo Yu. Multi-Class Data Description for Out-of-distribution Detection. *Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 1362–1370, aug 2020a. doi: 10.1145/3394486.3403189. URL <https://dl.acm.org/doi/10.1145/3394486.3403189>.
- Dongha Lee, Sehun Yu, and Hwanjo Yu. Multi-class data description for out-of-distribution detection. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, KDD '20, page 1362–1370, New York, NY, USA, 2020b. Association for Computing Machinery. ISBN 9781450379984. doi: 10.1145/3394486.3403189. URL <https://doi.org/10.1145/3394486.3403189>.
- Abhijit Bendale and Terrance E. Boult. Towards open set deep networks. *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, 2016-December:1563–1572, 12 2016. ISSN 10636919. doi: 10.1109/CVPR.2016.173.
- Yarin Gal and Zoubin Ghahramani. Dropout as a bayesian approximation: Representing model uncertainty in deep learning, 6 2016. ISSN 1938-7228. URL <https://proceedings.mlr.press/v48/gall16.html>.
- Shiyu Liang, Yixuan Li, and R. Srikant. Enhancing the reliability of out-of-distribution image detection in neural networks. *6th International Conference on Learning Representations, ICLR 2018 - Conference Track Proceedings*, 6 2017. URL <https://arxiv.org/abs/1706.02690v5>.
- Leo Grinsztajn, Edouard Oyallon, and Gael Varoquaux. Why do tree-based models still outperform deep learning on typical tabular data? *Advances in Neural Information Processing Systems*, 35:507–520, dec 2022.
- Achmad Ginanjar, Xue Li, and Wen Hua. Contrastive federated learning with tabular data silos, 2024. URL <https://arxiv.org/abs/2409.06123>.
- Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. pages 1597–1607, 11 2020. ISSN 2640-3498. URL <https://proceedings.mlr.press/v119/chen20j.html>.
- Talip Ucar, Ehsan Hajiramezanali, and Lindsay Edwards. Subtab: Subsetting features of tabular data for self-supervised representation learning. volume 23, 2021.
- Jonathan Huang, Vivek Rathod, Chen Sun, Menglong Zhu, Anoop Korattikara, Alireza Fathi, Ian Fischer, Zbigniew Wojna, Yang Song, Sergio Guadarrama, and Kevin Murphy. Speed/accuracy trade-offs for modern convolutional object detectors, 2017.

517 John Platt et al. Probabilistic outputs for support vector machines and comparisons to regularized likelihood
518 methods. *Advances in large margin classifiers*, 10(3):61–74, 1999.

519

520 Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q. Weinberger. On calibration of modern neural networks.
521 *34th International Conference on Machine Learning, ICML 2017*, 3:2130–2143, 6 2017. URL <https://arxiv.org/abs/1706.04599v2>.

522

523 Konstantin Kirchheim, Marco Filax, and Frank Ortmeier. Pytorch-ood: A library for out-of-distribution
524 detection based on pytorch. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern
525 Recognition (CVPR) Workshops*, pages 4351–4360, June 2022.

526

527 Dongha Lee, Sehun Yu, and Hwanjo Yu. Multi-class data description for out-of-distribution detection,
528 2020c. URL <https://github.com/donalee/DeepMCDD>.

529

529 Dennis W Ruck, Steven K Rogers, and Matthew Kabrisky. Feature selection using a multilayer perceptron.
530 *Journal of neural network computing*, 2(2):40–48, 1990.

531

531 Yury Gorishniy, Ivan Rubachev, Nikolay Kartashev, Daniil Shlenskii, Akim Kotelnikov, and Artem Babenko.
532 Tabr: Tabular deep learning meets nearest neighbors in 2023, 2023. URL <https://arxiv.org/abs/2307.14338>.

533

534 Günter Klambauer, Thomas Unterthiner, Andreas Mayr, and Sepp Hochreiter. Self-normalizing neural net-
535 works. *31st Conference on Neural Information Processing Systems*, 2017.

536

537 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser,
538 and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30,
539 2017.

540

541 Bin Li, Weihang Wei, Anselmo Ferreira, and Shunquan Tan. Rest-net: Diverse activation modules and
542 parallel subnets-based cnn for spatial image steganalysis. *IEEE Signal Processing Letters*, 25(5):650–
543 654, 2018.

544

544 Ruoxi Wang, Rakesh Shivanna, Derek Z. Cheng, Sagar Jain, Dong Lin, Lichan Hong, and Ed H. Chi. Dcn
545 v2: Improved deep and cross network and practical lessons for web-scale learning to rank systems. *The
546 Web Conference 2021 - Proceedings of the World Wide Web Conference, WWW 2021*, pages 1785–1797,
547 8 2020. doi: 10.1145/3442381.3450078. URL <http://arxiv.org/abs/2008.13535><http://dx.doi.org/10.1145/3442381.3450078>.

548

549 Weiping Song, Chence Shi, Zhiping Xiao, Zhijian Duan, Yewen Xu, Ming Zhang, and Jian Tang. Au-
550 toint: Automatic feature interaction learning via self-attentive neural networks. *International Confer-
551 ence on Information and Knowledge Management, Proceedings*, 10:1161–1170, 10 2018. doi: 10.1145/
552 3357384.3357925. URL <http://arxiv.org/abs/1810.11921><http://dx.doi.org/10.1145/3357384.3357925>.

553

554 Sergei Popov, Stanislav Morozov, and Artem Babenko. Neural oblivious decision ensembles for deep learn-
555 ing on tabular data. *8th International Conference on Learning Representations, ICLR 2020*, 9 2019. URL
556 <https://arxiv.org/abs/1909.06312v2>.

557

558 Sercan Arik and Tomas Pfister. Tabnet: Attentive interpretable tabular learning. *35th AAAI Conference on
559 Artificial Intelligence, AAAI 2021*, 8A:6679–6687, 8 2019. ISSN 2159-5399. doi: 10.1609/aaai.v35i8.
560 16826. URL <https://arxiv.org/abs/1908.07442v5>.

561

561 Sarkhan Badirli, Xuanqing Liu, Zhengming Xing, Avradeep Bhowmik, Khoa Doan, and Sathiya S. Keerthi.
562 Gradient boosting neural networks: Grownet. 2 2020. URL <https://arxiv.org/abs/2002.07971v2>.

563

564 Dara Bahri, Heinrich Jiang, Yi Tay, and Donald Metzler. SCARF: SELF-SUPERVISED CONTRASTIVE
565 LEARNING USING RANDOM FEATURE CORRUPTION. In *ICLR 2022 - 10th International Confer-*
566 *ence on Learning Representations*, 2022.

567

568 Barry Becker and Ronny Kohavi. Adult. UCI Machine Learning Repository, 1996. DOI:
569 <https://doi.org/10.24432/C5XW20>.

570

571

572 Isabelle Guyon, Lisheng Sun-Hosoya, Marc Boullé, Hugo Jair Escalante, Sergio Escalera, Zhengying Liu,
573 Damir Jajetic, Bisakha Ray, Mehreen Saeed, Michèle Sebag, Alexander Statnikov, Wei-Wei Tu, and Eve-
574 lyne Viegas. Analysis of the automl challenge series 2015-2018. In *AutoML*, Challenges in Machine
575 Learning. Springer, 2019.

576

577 Pierre Baldi, Peter Sadowski, and Daniel Whiteson. Searching for exotic particles in high-energy physics
578 with deep learning. *Nature Communications*, 5:4308, 2014. doi: 10.1038/ncomms5308.

579

580 Jan-Mark Geusebroek, Gertjan J. Burghouts, and Arnold W. M. Smeulders. The amsterdam library of
581 object images. *International Journal of Computer Vision*, 61(1):103–112, 2005. doi: 10.1023/B:VISI.
582 0000042993.50813.60.

583

584 PASCAL Challenge on Large Scale Learning. Epsilon Dataset: Simulated Physics Experiments. [http://](http://largescale.ml.tu-berlin.de/instructions/)
585 largescale.ml.tu-berlin.de/instructions/, 2008. Accessed: [Insert Access Date].

586

587

588 Jock A. Blackard and Denis J. Dean. Comparative accuracies of artificial neural networks and discrimi-
589 nant analysis in predicting forest cover types from cartographic variables. *Computers and Electronics in*
590 *Agriculture*, 24(3):131–151, 2000. doi: 10.1016/S0168-1699(99)00046-0.

591

592 R. Kelley Pace and Ronald Barry. Sparse spatial autoregressions. *Statistics & Probability Letters*, 33(3):
593 291–297, 1997.

594

595 Thierry Bertin-Mahieux, Daniel P.W. Ellis, Brian Whitman, and Paul Lamere. The million song dataset. In
596 *Proceedings of the 12th International Conference on Music Information Retrieval (ISMIR 2011)*, pages
597 591–596, Miami, Florida, USA, October 2011. URL [http://ismir2011.ismir.net/papers/](http://ismir2011.ismir.net/papers/OS6-1.pdf)
598 [OS6-1.pdf](http://ismir2011.ismir.net/papers/OS6-1.pdf).

599

600 Olivier Chapelle and Yi Chang. Yahoo! learning to rank challenge overview. In *Proceedings of the Learning*
601 *to Rank Challenge*, volume 14 of *Proceedings of Machine Learning Research*, pages 1–24. PMLR, 2011.
602 URL <http://proceedings.mlr.press/v14/chapelle11a.html>.

603

604

605 Tao Qin and Tie-Yan Liu. Introducing LETOR 4.0 datasets. *arXiv preprint arXiv:1306.2597*, 2013. URL
606 <https://arxiv.org/abs/1306.2597>.

607

608 anonymized. Tabular contrastive learning (tcl). [Online]. Available from: [https://github/](https://github.com/anonymized/tcl)
609 [anonymized](https://github.com/anonymized/tcl), July 12 2024.

610

611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657

A APPENDIX

Algorithm 1 OOD Detection and Train-Test Separation

Require: Dataset D , OOD detection methods $M = \{OpenMax, TemperatureScaling\}$
Ensure: In-distribution dataset D_{in} , Out-of-distribution dataset D_{ood}
for each method $m \in M$ **do**
 Apply m to D
 Obtain transformed data D_m
end for
for each D_m **do**
 Visualize histogram of D_m
 Manually set threshold t_m based on histogram
end for
Initialize $D_{in} \leftarrow \emptyset, D_{ood} \leftarrow \emptyset$
for each sample $x \in D$ **do**
 if $\forall m \in M : D_m(x) < t_m$ **then**
 $D_{in} \leftarrow D_{in} \cup \{x\}$
 else
 $D_{ood} \leftarrow D_{ood} \cup \{x\}$
 end if
end for
Validate OOD separation using linear regression
Compare model performance: $f(D_{in})$ vs $f(D_{ood})$
return D_{in}, D_{ood}

Algorithm 2 Tabular Contrastive Learning (TCL) Algorithm

Require: Tabular dataset \mathcal{D} , encoder E , decoder P , batch size N , temperature τ
for each epoch **do**
 for each batch $\mathcal{B} \in \mathcal{D}$ **do**
 for each sample $x_i \in \mathcal{B}$ **do**
 Create augmented views $x_i^1, x_i^2 = \text{Augment}(x_i)$
 end for
 $\mathcal{B}_{aug} = \{x_i^1, x_i^2 | i = 1, \dots, N\}$
 for each $\tilde{x}_i \in \mathcal{B}_{aug}$ **do**
 $x^e = E(\tilde{x}_i)$ ▷ Encode
 $x^p = P(x^e)$ ▷ Decode
 end for
 if inference **then return** x^e
 end if
 for $i = 1$ to $2N$ **do**
 $i^+ = (i + N) \bmod 2N$ ▷ Pair index
 $L_t = L_r + L_c + L_d$
 end for
 $L_{tN} = \frac{1}{N} \sum_1^N L_t$
 Update E and P by optimizing L_{tN}
 end for
 end for
end for
