# Revisiting Nearest Neighbor for Tabular Data: A Deep Tabular Baseline Two Decades Later

**Anonymous authors**
Paper under double-blind review

## Abstract

The widespread enthusiasm for deep learning has recently expanded into the domain of tabular data. Recognizing that the advancement in deep tabular methods is often inspired by classical methods, *e.g.*, integration of nearest neighbors into neural networks, we investigate whether these classical methods can be revitalized with modern techniques. We revisit a differentiable version of $K$-nearest neighbors (KNN) — Neighbourhood Components Analysis (NCA) — originally designed to learn a linear projection to capture semantic similarities between instances, and seek to gradually add modern deep learning techniques on top. Surprisingly, our implementation of NCA using SGD and without dimensionality reduction already achieves decent performance on tabular data, in contrast to the results of using existing toolboxes like scikit-learn. Further equipping NCA with deep representations and additional training stochasticity significantly enhances its capability, being on par with the leading tree-based method CatBoost and outperforming existing deep tabular models in both classification and regression tasks on 300 datasets. We conclude our paper by analyzing the factors behind these improvements, including loss functions, prediction strategies, and deep architectures.

## 1 Introduction

Tabular data, characterized by its structured format of rows and columns representing individual examples and features, is prevalent in domains like healthcare (Hassan et al., 2020) and e-commerce (Nederstigt et al., 2014). Motivated by the success of deep neural networks in fields like computer vision and natural language processing (Simonyan & Zisserman, 2015; Vaswani et al., 2017; Devlin et al., 2019), numerous deep models have been developed for tabular data to capture complex feature interactions (Cheng et al., 2016; Guo et al., 2017; Popov et al., 2020; Arik & Pfister, 2021; Gorishniy et al., 2021; Katzir et al., 2021; Chang et al., 2022; Chen et al., 2022; Hollmann et al., 2023).

Despite all these attempts, deep tabular models still struggle to match the accuracy of traditional machine learning methods like Gradient Boosting Decision Trees (GBDT) (Prokhorenkova et al., 2018; Chen & Guestrin, 2016) on tabular tasks. Such a fact raises our interest: *to excel in tabular tasks, perhaps deep methods could draw inspiration from traditional methods.* Indeed, several deep tabular methods have demonstrated promising results along this route. Gorishniy et al. (2021); Kadra et al. (2021) consulted classical tabular techniques to design specific MLP architectures and weight regularization strategies, significantly boosting MLPs' accuracy on tabular datasets. Recently, inspired by non-parametric methods (Mohri et al., 2012), TabR (Gorishniy et al., 2024) retrieves neighbors from the entire training set and constructs instance-specific scores with a Transformer-like architecture, leveraging relationships between instances for tabular predictions.

*We follow this route but from a different direction.* Instead of incorporating classic techniques into the already complex deep models, we perform an Occam's-razor-style exploration — starting from the classic method and gradually increasing its complexity by adding modern deep techniques. We hope such an exploration could reveal the key components from both worlds to excel in tabular tasks.

To this end, we build upon TabR (Gorishniy et al., 2024) and choose to start from a classical, differentiable version of $K$-nearest neighbors (KNN) named Neighbourhood Component Analysis (NCA) (Goldberger et al., 2004). NCA optimizes the KNN prediction accuracy of a target instance by learning a linear projection, ensuring that semantically similar instances are closer than dissimilar ones. Its differentiable nature makes it a suitable backbone for adding deep learning modules.
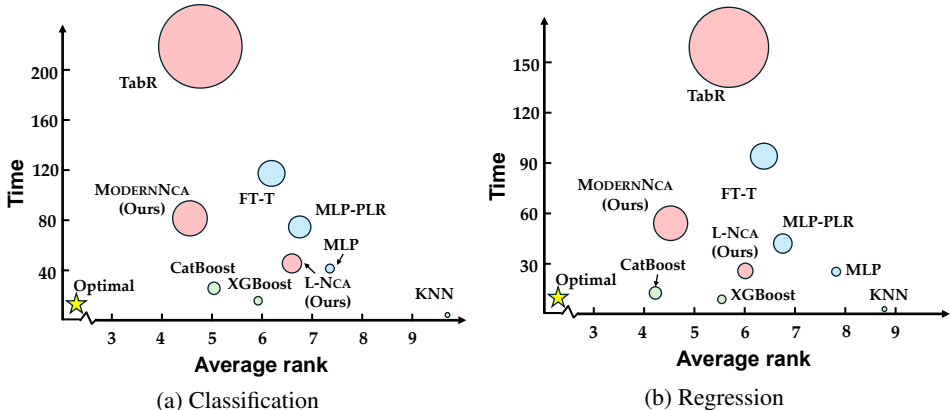
Figure 1: Performance-Efficiency-Memory comparison between MODERNNCA and existing methods on classification (a) and regression (b) datasets. Representative tabular prediction methods, including the *classical* methods (in green), the *parametric* deep methods (in blue), and the *non-parametric/neighborhood-based* deep methods (in red), are investigated, based on their records over 300 datasets in Table 1 and Figure 2. The average rank among these eight methods is used as the performance measure. We calculate the average training time (in seconds) and the memory usage of the model (denoted by the radius of the circles, where the larger the circle, the bigger the model). MODERNNCA achieves high training speed compared to other deep tabular models and has a relatively lower memory usage. L-NCA is our *improved linear* version of NCA.

Our first attempt is to re-implement NCA, using deep learning libraries like PyTorch (Paszke et al., 2019). Interestingly, by replacing the default L-BFGS optimizer (Liu & Nocedal, 1989) in scikit-learn (Pedregosa et al., 2011)[1] with stochastic gradient descent (SGD), we already witnessed a notable accuracy boost on tabular tasks. Further enabling NCA to learn a linear projection into a larger dimensionality (hence not dimensionality reduction) and use a soft nearest neighbor inference rule (Salakhutdinov & Hinton, 2007; Frosst et al., 2019a) bring another gain, making NCA on par with deep methods like MLP. (See section 6 for detailed ablation studies and discussions.)

Our second attempt is to replace the linear projection with a neural network for nonlinear embeddings. As NCA's objective function involves the relationship of an instance to all the other training instances, a naive implementation would incur a huge computational burden. We thus employ a stochastic neighborhood sampling (SNS) strategy, randomly selecting a subset of training data as candidate neighbors in each mini-batch. We show that SNS not only improves training efficiency but enhances the model's generalizability, as it introduces additional stochasticity (beyond SGD) in training.

Putting things together, along with the use of a pre-defined feature transform on numerical tabular entries (Gorishniy et al., 2022), our deep NCA implementation, MODERNNCA, achieves remarkably encouraging empirical results. Evaluated on 300 tabular datasets, MODERNNCA is ranked first in classification tasks and just shy of CatBoost (Prokhorenkova et al., 2018) in regression tasks while outperforming other tree-based and deep tabular models. Figure 1 further shows that MODERNNCA well balances training efficiency (with lower training time compared to other deep tabular models), generalizability (with higher average accuracy), and memory efficiency. We also provide a detailed ablation study and discussion on MODERNNCA, comparing different loss functions, training and prediction strategies, and deep architectures, aiming to systematically reveal the impacts of deep learning techniques on NCA, after its release in 2004. In sum, our contributions are two-folded:

- We revisit the classical nearest neighbor approach NCA and systematically explore ways to improve it using modern deep learning techniques.
- Our proposed MODERNNCA achieves outstanding performance in both classification and regression tasks, essentially serving as a strong deep baseline for tabular tasks.

**Remark.** In conducting this study, we become aware of several prior attempts to integrate neural networks into NCA (Salakhutdinov & Hinton, 2007; Min et al., 2010). However, their results and applicability were downplayed by tree-based methods, and we attribute this to the less powerful deep-learning techniques two decades ago (*e.g.*, restricted Boltzmann machine). In other words, our work can be viewed as a revisit of these attempts from the lens of modern deep-learning techniques.

---

[1]We note that the original NCA paper (Goldberger et al., 2004) did not specify the optimizer.

While our study is largely *empirical*, it should not be seen as a weakness. For years, nearest-neighbor-based methods (though with solid theoretical foundations) have been overlooked in tabular data, primarily due to their low competitiveness with tree-based methods. We hope that our thorough exploration of deep learning techniques for nearest neighbors and the outcome — a strong tabular baseline on par with the leading CatBoost (Prokhorenkova et al., 2018) — would revitalize nearest neighbors and open up new research directions, ideally theoretical foundations behind the improvements.

## 2 RELATED WORK

**Learning with Tabular Data**. Tabular data is a common format across various applications such as click-through rate prediction (Richardson et al., 2007) and time-series forecasting (Ahmed et al., 2010). Tree-based methods like XGBoost (Chen & Guestrin, 2016), LightGBM (Ke et al., 2017), and CatBoost (Prokhorenkova et al., 2018) have proven effective at capturing feature interactions and are widely used in real-world applications. Recognizing the ability of deep neural networks to learn feature representations from raw data and make nonlinear predictions, recent methods have applied deep learning techniques to tabular models (Cheng et al., 2016; Guo et al., 2017; Popov et al., 2020; Borisov et al., 2022; Arik & Pfister, 2021; Kadra et al., 2021; Katzir et al., 2021; Chen et al., 2022). For instance, deep architectures such as residual networks and transformers have been adapted for tabular prediction (Gorishniy et al., 2021; Hollmann et al., 2023). Moreover, data augmentation strategies have been introduced to mitigate overfitting in deep models (Ucar et al., 2021; Bahri et al., 2022; Rubachev et al., 2022). Deep tabular models have demonstrated competitive performance across a wide range of applications. However, researchers have observed that deep models still face challenges in capturing high-order feature interactions as effectively as tree-based models (Grinsztajn et al., 2022; McElfresh et al., 2023; Ye et al., 2024a).

**NCA Variants**. Nearest Neighbor approaches make predictions based on the relationships between an instance and its neighbors in the training set. Instead of identifying neighbors using raw features, NCA employs a differentiable Nearest Neighbor loss function (also known as soft-NN loss) to learn a linear projection for better distance measurement (Goldberger et al., 2004). Several works have extended this idea with alternative loss functions (Globerson & Roweis, 2005; Tarlow et al., 2013), while others explore NCA variants for data visualization (Venna et al., 2010). A few nonlinear extensions of NCA, developed over a decade ago, demonstrated a bit improved performance on image classification tasks using architecture like restricted Boltzmann machines (Salakhutdinov & Hinton, 2007; Min et al., 2010). For visual tasks, the entanglement effects of soft-NN loss on deep learned representations have been analyzed (Frosst et al., 2019b), and variants of this loss have been applied to few-shot learning scenarios (Vinyals et al., 2016; Laenen & Bertinetto, 2021). The effectiveness of NCA variants in fields like image recognition suggests untapped potential, motivating our revisit of this method with modern deep learning techniques for tabular data.

**Metric Learning**. NCA is a form of metric learning (Xing et al., 2002), where a projection is learned to pull similar instances closer together and push dissimilar ones farther apart, leading to improved classification and regression performance with KNN (Davis et al., 2007; Weinberger & Saul, 2009; Kulis, 2013; Bellet et al., 2015). Initially applied to tabular data, metric learning has evolved into a valuable tool, particularly when integrated with deep learning techniques, across domains like image recognition (Schroff et al., 2015; Sohn, 2016; Song et al., 2016; Khosla et al., 2020), person re-identification (Yi et al., 2014; Yang et al., 2018), and recommendation systems (Hsieh et al., 2017; Wei et al., 2023). Recently, TabR (Gorishniy et al., 2024) introduced a variant of transformer architecture to retrieve neighbors for each instance, enhancing tabular prediction tasks. Despite its promising results, the high computational cost of neighborhood selection and the complexity of its architecture limit the practicality of TabR. In contrast, our paper revisits NCA and proposes a simpler deep tabular baseline that maintains efficient training speeds without sacrificing performance.

## 3 PRELIMINARY

In this section, we first introduce the task learning with tabular data. We then provide a brief overview of NCA (Goldberger et al., 2004) and TabR (Gorishniy et al., 2024).

### 3.1 LEARNING WITH TABULAR DATA

A labeled tabular dataset is formatted as $N$ examples (rows in the table) and $d$ features/attributes (columns in the table). An instance $\boldsymbol{x}_i$ is depicted by its $d$ feature values. There are two kinds of features: the numerical (continuous) ones and categorical (discrete) ones. Given $x_{i,j}$ as the $j$-th feature of instance $\boldsymbol{x}_i$, we use $x_{i,j}^{\text{num}} \in \mathbb{R}$ and $\boldsymbol{x}_{i,j}^{\text{cat}}$ to denote numerical (*e.g.*, the height of a person) and categorical (*e.g.*, the gender of a person) feature values of an instance, respectively. The categorical features are usually transformed in a one-hot manner, *i.e.*, $\boldsymbol{x}_{i,j}^{\text{cat}} \in \{0,1\}^{K_j}$, where the index of value 1 indicates the category among the $K_j$ options. We assume the instance $\boldsymbol{x}_i \in \mathbb{R}^d$ w.l.o.g. and will explore other encoding strategies later. Each instance is associated with a label $y_i$, where $y_i \in [C] = \{1, \ldots, C\}$ in a multi-class classification task and $y_i \in \mathbb{R}$ in a regression task.

Given a tabular dataset $\mathcal{D} = \{(\boldsymbol{x}_i, y_i)\}_{i=1}^N$, we aim to learn a model $f$ on $\mathcal{D}$ that maps $\boldsymbol{x}_i$ to its label $y_i$. We measure the quality of $f$ by the joint likelihood over $\mathcal{D}$, *i.e.*, $\max_f \prod_{(\boldsymbol{x}_i, y_i) \in \mathcal{D}} \Pr(y_i \mid f(\boldsymbol{x}_i))$. The objective could be reformulated in the form of negative log-likelihood of the true labels,

$$\min_f \sum_{(\boldsymbol{x}_i, y_i) \in \mathcal{D}} - \log \Pr(y_i \mid f(\boldsymbol{x}_i)) = \sum_{(\boldsymbol{x}_i, y_i) \in \mathcal{D}} \ell(y_i, \; \hat{y}_i = f(\boldsymbol{x}_i)) \;, \tag{1}$$

or equivalently, the discrepancy between the predicted label $\hat{y}_i$ and the true label $y_i$ measured by the loss $\ell(\cdot, \cdot)$, *e.g.*, cross-entropy. We expect the learned model $f$ is able to extend its ability to unseen instances sampled from the same distribution as $\mathcal{D}$. $f$ could be implemented with classical methods such as SVM and tree-based approaches or MLPs.

### 3.2 NEAREST NEIGHBOR FOR TABULAR DATA

**KNN** is one of the most representative non-parametric tabular models for classification and regression — making predictions based on the labels of the nearest neighbors (Bishop, 2006; Mohri et al., 2012). In other words, the prediction $f(\boldsymbol{x}_i; \mathcal{D})$ of the model $f$ conditions on the whole training set. Given an instance $\boldsymbol{x}_i$, KNN calculates the distance between $\boldsymbol{x}_i$ and other instances in $\mathcal{D}$. Assume the $K$ nearest neighbors are $\mathcal{N}(\boldsymbol{x}_i; \mathcal{D}) = \{(\boldsymbol{x}_1, y_1), \ldots, (\boldsymbol{x}_K, y_K)\}$, then, the label $y_i$ of $\boldsymbol{x}_i$ is predicted based on those labels in the neighbor set $\mathcal{N}(\boldsymbol{x}_i; \mathcal{D})$. For classification task $\hat{y}_i$ is the majority voting of labels in $\mathcal{N}(\boldsymbol{x}_i; \mathcal{D})$ while is the average of those labels in regression tasks.

The distance $\text{dist}(\boldsymbol{x}_i, \boldsymbol{x}_j)$ in KNN determines the set of nearest neighbors $\mathcal{N}(\boldsymbol{x}_i; \mathcal{D})$, which is one of its key factors. The Euclidean distance between a pair $(\boldsymbol{x}_i, \boldsymbol{x}_j)$ is $\text{dist}(\boldsymbol{x}_i, \boldsymbol{x}_j) = \sqrt{(\boldsymbol{x}_i - \boldsymbol{x}_j)^\top (\boldsymbol{x}_i - \boldsymbol{x}_j)}$. A distance metric that reveals the characteristics of the dataset will improve KNN and lead to more accurate predictions (Xing et al., 2002; Davis et al., 2007; Weinberger & Saul, 2009; Bellet et al., 2015).

**Neighbourhood Component Analysis (NCA).** NCA focuses on the classification task (Goldberger et al., 2004). According to the 1NN rule, NCA defines the probability that $\boldsymbol{x}_j$ locates in the neighborhood of $\boldsymbol{x}_i$ by

$$\Pr(\boldsymbol{x}_j \in \mathcal{N}(\boldsymbol{x}_i; \mathcal{D}) \mid \boldsymbol{x}_i, \mathcal{D}, \boldsymbol{L}) = \frac{\exp\left(-\text{dist}^2(\boldsymbol{L}^\top \boldsymbol{x}_i, \; \boldsymbol{L}^\top \boldsymbol{x}_j)\right)}{\sum_{(\boldsymbol{x}_l, y_l) \in \mathcal{D}, \boldsymbol{x}_l \neq \boldsymbol{x}_i} \exp\left(-\text{dist}^2(\boldsymbol{L}^\top \boldsymbol{x}_i, \; \boldsymbol{L}^\top \boldsymbol{x}_l)\right)} \;. \tag{2}$$

Then, the posterior probability that an instance $\boldsymbol{x}_i$ is classified as the class $y_i$ is:

$$\Pr(\hat{y}_i = y_i \mid \boldsymbol{x}_i, \mathcal{D}, \boldsymbol{L}) = \sum_{(\boldsymbol{x}_j, y_j) \in \mathcal{D} \wedge y_j = y_i} \Pr(\boldsymbol{x}_j \in \mathcal{N}(\boldsymbol{x}_i; \mathcal{D}) \mid \boldsymbol{x}_i, \mathcal{D}, \boldsymbol{L}) \;. \tag{3}$$

$\boldsymbol{L} \in \mathbb{R}^{d \times d'}$ is a linear projection usually with $d' \leq d$, which reduces the dimension of the raw input. Therefore, the posterior that an instance $\boldsymbol{x}_i$ belongs to the class $y_i$ depends on its similarity (measured by the negative squared Euclidean distance in the space projected by $\boldsymbol{L}$) between its neighbors from class $y_i$ in $\mathcal{D}$. Equation 3 approximates the expected leave-one-out error for $\boldsymbol{x}_i$, and the original NCA maximizes the **sum of** $\Pr(\hat{y}_i = y_i \mid \boldsymbol{x}_i, \mathcal{D}, \boldsymbol{L})$ over all instances in $\mathcal{D}$. Instead of considering all instances in the neighborhood equally, this objective mimics a soft version of KNN, where all instances in the training set are weighted (nearer neighbors have more weight) for the nearest neighbor decision. In the test stage, KNN is applied to classify an unseen instance in the space projected by $\boldsymbol{L}$.

**TabR** is a deep tabular method where the neighbors of an instance $\boldsymbol{x}_i$ are retrieved with deep neural networks. In detail, TabR defines the contribution of $(\boldsymbol{x}_j, y_j)$ to the predicted label $\hat{y}_i$ of $\boldsymbol{x}_i$ as

$$s(\boldsymbol{x}_i, \boldsymbol{x}_j, y_j) = \boldsymbol{W}\boldsymbol{y}_j + \mathrm{T}(\boldsymbol{L}^\top \boldsymbol{x}_j - \boldsymbol{L}^\top \boldsymbol{x}_i) \ . \tag{4}$$

The transformation T is the composition of a linear layer without bias, dropout, ReLU, and another linear layer. $\boldsymbol{W}$ is a linear projection and $\boldsymbol{y}_j$ is the encoded label vector of $y_j$. The instance-specific scores are then weighted to obtain $\hat{y}_i = \sum_{(\boldsymbol{x}_j, y_j) \in \mathcal{D}} \alpha_j \cdot s(\boldsymbol{x}_i, \boldsymbol{x}_j, y_j)$. The weight $\alpha_j \propto \{-\operatorname{dist}(\boldsymbol{L}^\top \boldsymbol{x}_j, \ \boldsymbol{L}^\top \boldsymbol{x}_i)\}$, and is normalized by a softmax operator. Please refer to Gorishniy et al. (2024) for additional details such as the layer normalization over instances, encoding of numerical attributes, and the selection of K nearest neighbors in the summation.

# 4 MODERNNCA

Given the promising results of TabR on tabular data, we take the original NCA as our starting point and gradually enhance its complexity by incorporating modern deep learning techniques. This Occam's-razor-style exploration may allow us to identify the key components that lead to strong performance in tabular tasks, drawing insights from both classical and deep tabular models. In the following, we introduce our proposed MODERNNCA (abbreviated as M-NCA) through two key attempts to improve upon the original NCA.

## 4.1 THE FIRST ATTEMPT

We generalize the projection in Equation 2 by introducing a transformation $\phi$, which maps $\boldsymbol{x}_i$ into a space with dimensionality $d'$. To remain consistent with the original NCA, we initially define $\phi$ as a linear layer, *i.e.*, $\phi(\boldsymbol{x}_i) = \mathrm{Linear}(\boldsymbol{x}_i)$, consisting of a linear projection and a bias term.

**Learning Objective**. Assume the label $y_j$ is continuous in regression tasks and in one-hot form for classification tasks. We modify Equation 3 as follows:

$$\hat{y}_i = \sum_{(\boldsymbol{x}_j, y_j) \in \mathcal{D}} \frac{\exp\left(-\operatorname{dist}^2(\phi(\boldsymbol{x}_i), \ \phi(\boldsymbol{x}_j))\right)}{\sum_{(\boldsymbol{x}_l, y_l) \in \mathcal{D}, \boldsymbol{x}_l \neq \boldsymbol{x}_i} \exp\left(-\operatorname{dist}^2(\phi(\boldsymbol{x}_i), \ \phi(\boldsymbol{x}_l))\right)} y_j \ . \tag{5}$$

This formulation ensures that similar instances (based on their distance in the embedding space mapped by $\phi$) yield closer predictions. For classification, Equation 5 generalizes Equation 3, predicting the label of a target instance by computing a weighted average of its neighbors across the $C$ classes. Here, $\hat{y}_i \in \mathbb{R}^C$ is a probability vector representing $\{\Pr(\hat{y}_i = c \mid \boldsymbol{x}_i, \mathcal{D}, \phi)\}_{c \in [C]}$. In regression tasks, the prediction is the weighted sum of scalar labels from the neighborhood.

By combining Equation 3 with Equation 1, we define $\ell$ in Equation 1 as negative log-likelihood for classification and mean square error for regression. This classification loss is also known as the soft Nearest Neighbor (soft-NN) loss (Frosst et al., 2019b; Khosla et al., 2020) for visual tasks. Different from Goldberger et al. (2004); Salakhutdinov & Hinton (2007) that used **sum of probability** as in the original NCA's loss, we find **sum of log probability** provides better performance on tabular data.

**Prediction Strategy**. For a test instance, the original NCA projects all instances using the learned $\phi$ and applies KNN to classify the test instance based on its neighbors from the entire training set $\mathcal{D}$. Instead of employing the traditional "hard" KNN approach, we adopt the soft-NN rule (Equation 5) to estimate the label posterior, applicable to both classification and regression. Specifically, in the classification case, Equation 5 produces a $C$-dimensional vector, with the index of the maximum value indicating the predicted class. For regression, $\hat{y}_i$ directly corresponds to the predicted value.

Furthermore, we do not limit the mapping to dimensionality reduction. The linear projection $\phi$ can transform $\boldsymbol{x}_i$ into a higher-dimensional space if necessary. We also replace the L-BFGS optimizer (used in scikit-learn) with stochastic gradient descent (SGD) for better scalability and performance.

These modifications result in a notable accuracy boost for NCA on tabular tasks, making it competitive with deep models like MLP. We refer to this improved version of (linear) NCA as L-NCA.

## 4.2 THE SECOND ATTEMPT

We further enhance L-NCA by incorporating modern deep learning techniques, leading to our strong deep tabular baseline, MODERNNCA (M-NCA).

**Architectures.** To introduce nonlinearity into the model, we first enhance the transformation $\phi$ in subsection 4.1 with multiple nonlinear layers appended. Specifically, we define a one-layer nonlinear mapping as a sequence of operators following Gorishniy et al. (2021), consisting of one-dimensional batch normalization (Ioffe & Szegedy, 2015), a linear layer, ReLU activation, dropout (Srivastava et al., 2014), and another linear layer. In other words, the input $\boldsymbol{x}_i$ will be transformed by

$$g(\boldsymbol{x}_i) = \text{Linear}\left(\text{Dropout}\left(\left(\text{ReLU}\left(\text{Linear}\left(\text{BatchNorm}\left(\boldsymbol{x}_i\right)\right)\right)\right)\right)\right) . \tag{6}$$

One or more layers of such a block $g$ can be appended on top of the original linear layer in subsection 4.1 to implement the final nonlinear mapping $\phi$, which further incorporates an additional batch normalization at the end to calibrate the output embedding. Empirical results show that batch normalization outperforms other normalization strategies, such as layer normalization (Ba et al., 2016), in learning a robust latent embedding space.

For categorical input features, we use one-hot encoding, and for numerical features, we leverage PLR (lite) encoding, following TabR (Gorishniy et al., 2024). PLR encoding combines periodic embeddings, a linear layer, and ReLU to project instances into a high-dimensional space, thereby increasing the model's capacity with additional nonlinearity (Gorishniy et al., 2022). PLR (lite) restricts the linear layer to be shared across all features, balancing complexity and efficiency.

**Stochastic Neighborhood Sampling**. SGD is commonly applied to optimize deep neural networks — a mini-batch of instances is sampled, and the average instance-wise loss in the mini-batch is calculated for back-propagation. However, the instance-wise loss based on the predicted label in Equation 5 involves pairwise distances between an instance in the mini-batch and the entire training set $\mathcal{D}$, imposing a significant computational burden.

To accelerate the training speed of MODERNNCA, we propose a Stochastic Neighborhood Sampling (SNS) strategy. In SNS, a subset $\hat{\mathcal{D}}$ of the training set $\mathcal{D}$ is randomly sampled for each mini-batch, and only distances between instances in the mini-batch and this subset are calculated. In other words, $\hat{\mathcal{D}}$ replaces $\mathcal{D}$ in Equation 5, and only the labels in $\hat{\mathcal{D}}$ are used to predict the label of a given instance during training. During inference, however, the model resumes the searches for neighbors using the entire training set $\mathcal{D}$. Unlike deep metric learning methods that only consider pairs of instances within a sampled mini-batch (Schroff et al., 2015; Song et al., 2016; Sohn, 2016), *i.e.*, $\hat{\mathcal{D}}$ is the mini-batch, our SNS approach retains both efficiency and diversity in the selection of neighbor candidates.

We empirically observed that SNS not only increases the training efficiency of MODERNNCA, since fewer examples are utilized for back-propagation, but also improves the generalization ability of the learned mapping $\phi$. We attribute the improvement to the fact that $\phi$ is learned on more difficult, stochastic prediction tasks. The resulting $\phi$ thus becomes more robust to the potentially noisy and unstable neighborhoods in the test scenario. The influence of sampling ratio and other sampling strategies are investigated in detail in the experiments.

**Distance Function.** Empirically, we find that using the Euclidean distance instead of its squared form in Equation 5 leads to further performance improvements. Therefore, we adopt Euclidean distance as the default. Comparisons of various distance functions are provided in the appendix.

## 5 EXPERIMENTS

### 5.1 SETUPS

**Datasets.** We validate MODERNNCA over 300 datasets from a recently released large-scale tabular benchmark (Ye et al., 2024a), which includes 120 classification datasets and 180 regression datasets collected from UCI, OpenML (Vanschoren et al., 2014), Kaggle, and other sources.

**Evaluation.** We follow the evaluation protocol from Gorishniy et al. (2021; 2024). Each dataset is randomly split into training, validation, and test sets in proportions of 64%/16%/20%, respectively. For each dataset, we train each model using 15 different random seeds and calculate the average

Figure 2: The critical difference diagrams based on the Wilcoxon-Holm test with a significance level of 0.05 to detect pairwise significance for both classification tasks (evaluated using accuracy) and regression tasks (evaluated using RMSE).

Table 1: The Win/Tie/Lose ratio between MODERNNCA and 20 comparison methods across the 300 datasets, covering both classification (based on accuracy) and regression tasks (based on RMSE). This ratio is determined using a significant $t$-test at a 95% confidence interval.

| Method | Win | Tie | Lose | Method | Win | Tie | Lose |
|---|---|---|---|---|---|---|---|
| SVM | 0.78 | 0.13 | 0.10 | KNN | 0.79 | 0.07 | 0.14 |
| SwitchTab (Wu et al., 2024) | 0.88 | 0.09 | 0.03 | DANets (Chen et al., 2022) | 0.74 | 0.18 | 0.08 |
| NODE (Popov et al., 2020) | 0.70 | 0.15 | 0.15 | Tangos (Jeffares et al., 2023) | 0.66 | 0.20 | 0.14 |
| TabCaps (Chen et al., 2023) | 0.64 | 0.23 | 0.13 | PTaRL (Ye et al., 2024b) | 0.62 | 0.22 | 0.16 |
| DCNv2 (Wang et al., 2021) | 0.62 | 0.20 | 0.18 | MLP (Gorishniy et al., 2021) | 0.61 | 0.23 | 0.15 |
| ResNet (Gorishniy et al., 2021) | 0.59 | 0.30 | 0.11 | MLP-PLR (Gorishniy et al., 2022) | 0.57 | 0.27 | 0.16 |
| RandomForest | 0.57 | 0.18 | 0.26 | ExcelFormer (Chen et al., 2024) | 0.56 | 0.28 | 0.16 |
| SAINT (Somepalli et al., 2022) | 0.56 | 0.26 | 0.19 | FT-T (Gorishniy et al., 2021) | 0.50 | 0.28 | 0.23 |
| XGBoost (Chen & Guestrin, 2016) | 0.49 | 0.19 | 0.32 | LightGBM Ke et al. (2017) | 0.46 | 0.21 | 0.33 |
| TabR (Gorishniy et al., 2024) | 0.42 | 0.34 | 0.24 | CatBoost (Prokhorenkova et al., 2018) | 0.38 | 0.23 | 0.39 |

performance on the test set. For classification tasks, we consider accuracy (higher is better), and for regression tasks, we use Root Mean Square Error (RMSE) (lower is better). To summarize overall model performance, we report the average performance rank across all methods and datasets (lower ranks are better), following Delgado et al. (2014); McElfresh et al. (2023). Additionally, we conduct statistical $t$-tests to determine whether the differences between MODERNNCA and other methods are statistically significant.

**Comparison Methods.** We compare MODERNNCA with 20 approaches among three different categories. (For brevity, only 8 of them are shown in Figure 1.) First, we consider **classical parametric methods**, including linear SVM and tree-based methods like RandomForest, XGBoost (Chen & Guestrin, 2016), LightGBM Ke et al. (2017), and CatBoost (Prokhorenkova et al., 2018). Then, we consider **parametric deep models**, including NODE (Popov et al., 2020), MLP (Kadra et al., 2021; Gorishniy et al., 2021), ResNet (Gorishniy et al., 2021), SAINT (Somepalli et al., 2022), DCNv2 (Wang et al., 2021), FT-Transformer (Gorishniy et al., 2021), DANets (Chen et al., 2022), MLP-PLR (Gorishniy et al., 2022), TabCaps (Chen et al., 2023), Tangos (Jeffares et al., 2023), PTaRL (Ye et al., 2024b), SwitchTab (Wu et al., 2024), and ExcelFormer (Chen et al., 2024). For **neighborhood-based methods**, we consider KNN and TabR (Gorishniy et al., 2024). For a fair comparison, the same PLR numerical encoding is applied in MLP-PLR, TabR, and MODERNNCA.

**Implementation Details.** We pre-process all datasets following Gorishniy et al. (2021). For all deep methods, we set the batch size as 1024. The hyper-parameters of all methods are searched based on the training and validation set via Optuna (Akiba et al., 2019) following Gorishniy et al. (2021; 2024) over 100 trials. We set the ranges of the hyper-parameters for the compared methods following Gorishniy et al. (2021; 2024) and their official codes. The best-performed hyper-parameters are fixed during the final 15 seeds. Since the sampling rate of SNS effectively enhances the performance and reduces the training overhead, we treat it as a hyper-parameter and search within the range of [0.05, 0.6].

## 5.2 MAIN RESULTS

The comparison results between MODERNNCA, L-NCA, and six representative methods are presented in Figure 1. All methods are evaluated across three aspects: performance (average performance rank), average training time, and average memory usage across all datasets. While some models,

such as TabR, exhibit strong performance, they require significantly longer training times. In contrast, MODERNNCA strikes an excellent balance across various evaluation criteria.

We also applied the Wilcoxon-Holm test (Demsar, 2006) to assess pairwise significance among all methods for both classification and regression tasks. The results are shown in Figure 2. For classification tasks (shown in the left part of Figure 2), MODERNNCA consistently outperforms tree-based methods like XGBoost in most cases, demonstrating that its deep neural network architecture is more effective at capturing nonlinear relationships. Furthermore, compared to deep tabular models such as FT-T and MLP-PLR, MODERNNCA maintains its superiority. When combined with the results in Figure 1, these observations validate the effectiveness of MODERNNCA. It achieves performance on par with the leading tree-based method, CatBoost, while outperforming existing deep tabular models in both classification and regression tasks across 300 datasets.

Additionally, we calculated the Win/Tie/Lose ratio between MODERNNCA and other comparison methods across the 300 datasets. If two methods show no significant difference (based on a $t$-test at a 95% confidence interval), they are considered tied. Otherwise, one method is declared the winner based on the comparison of their average performance. Given the no free lunch theorem, it is challenging for any single method to statistically outperform others across all cases. Nevertheless, MODERNNCA demonstrates superior performance in most cases. For instance, MODERNNCA outperforms TabR on 126 datasets, ties on 102 datasets, and does so with a simpler architecture and shorter training time. Compared to CatBoost, MODERNNCA wins on 114 datasets and ties on 69 datasets. These results indicate that MODERNNCA serves as an effective and competitive deep learning baseline for tabular data.

## 6 ANALYSES AND ABLATION STUDIES OF MODERNNCA

In this section, we analyze the sources of improvement in MODERNNCA. All experiments are conducted on a tiny tabular benchmark comprising 45 datasets, as introduced in (Ye et al., 2024a). The benchmark consists of 27 classification datasets and 18 regression datasets. The average rank of various tabular methods on this benchmark closely aligns with the results observed on the larger set of 300 datasets, as detailed in (Ye et al., 2024a).

### 6.1 IMPROVEMENTS FROM NCA TO L-NCA

We begin with the original NCA (Goldberger et al., 2004), using the scikit-learn implementation (Pedregosa et al., 2011). We progressively replace key components in NCA and assess the resulting performance improvements. Since the original NCA only targets classification tasks, this subsection focuses on the 27 classification datasets in the tiny benchmark. To ensure a fair comparison, we re-implement the original NCA using the deep learning framework PyTorch (Paszke et al., 2019), denoting this baseline version as "NCAv0".

**Does Projection to a Higher Dimension Help?** In the scikit-learn implementation, NCA is constrained to perform dimensionality reduction, *i.e.*, $d' \leq d$ for the projection $\boldsymbol{L}$. We remove this constraint, allowing NCA to project into higher dimensions, and refer to this version as "NCAv1". Although higher dimensions by linear projections do not inherently enhance the representation ability of the squared Euclidean distance, the improvements in average performance rank from NCAv0 to NCAv1 (shown in Table 2) indicate that projecting to a higher dimension facilitates the optimization of this non-convex problem and improves generalization.

**Does Stochastic Gradient Descent Help?** Stochastic gradient descent (SGD) is a widely used optimizer in deep learning. To explore whether SGD can improve NCA's performance, we replace the default L-BFGS optimizer used in scikit-learn with SGD (without momentum) and denote this variant as "NCAv2". The performance improvements from NCAv1 to NCAv2 in Table 2 indicate that SGD makes NCA more effective in tabular data tasks.

**The Influence of the Loss Function.** The original NCA maximizes the expected leave-one-out accuracy as shown in Equation 3. In contrast, we minimize the negative log version of this objective as described in Equation 1. Although the log version for classification tasks was mentioned in Goldberger et al. (2004); Salakhutdinov & Hinton (2007), the original NCA preferred the leave-one-out formulation for better performance. We denote the variant with the modified loss function as

Table 2: Comparison of the average rank of (the linear) NCA variants and (the nonlinear) MLP across 27 classification datasets in the tiny-benchmark. The check marks indicate the differences in components among the variants. The average rank represents the overall performance of a method across all datasets, with lower ranks indicating better performance. The final variant, NCAv4, corresponds to the L-NCA version discussed in our paper.

| | High dimension | SGD optimizer | Log loss | Soft-NN prediction | Average rank |
|---|---|---|---|---|---|
| NCAv0 | | | | | 4.400 |
| NCAv1 | ✓ | | | | 3.708 |
| NCAv2 | ✓ | ✓ | | | 3.296 |
| NCAv3 | ✓ | ✓ | ✓ | | 3.192 |
| NCAv4 | ✓ | ✓ | ✓ | ✓ | 2.962 |
| MLP | ✓ | ✓ | ✓ | | 3.000 |

Table 3: Comparison among various configurations of the deep architectures used to implement $\phi$, where MLP is the default choice in MODERNNCA. We show the change in average performance rank (lower is better) across the four configurations on the 45 datasets in the tiny benchmark.

| | MLP | Linear | w/ LayerNorm | ResNet |
|---|---|---|---|---|
| Classification | 2.333 | 2.778 | 2.537 | 2.352 |
| Regression | 2.333 | 2.433 | 2.528 | 2.806 |

"NCAv3". As shown in Table 2 (NCAv2 vs. NCAv3), we find that using the log version slightly improves performance, especially when combined with deep architectures used in MODERNNCA. Further comparisons with alternative objectives are provided in the appendix.

**The Influence of the Prediction Strategy.** During testing, rather than applying a "hard" KNN with the learned embeddings as in standard metric learning, we adopt a soft nearest neighbor (soft-NN) inference rule, consistent with the training phase. This variant, using soft-NN for prediction, is referred to as "NCAv4", which is equivalent to the "L-NCA" version defined in subsection 4.1. Based on the change of average performance rank in Table 2, this modified prediction strategy further enhances NCA's classification performance, bringing linear NCA surpassing deep models like MLP.

## 6.2 IMPROVEMENTS FROM L-NCA TO M-NCA

In this subsection, we investigate the influence of architectures and encoding strategies to systematically reveal the impacts of more deep learning techniques on NCA.

**Linear vs. Deep Architectures**. We first investigate the architecture design for $\phi$ in MODERNNCA, where one or more layers of blocks $g(\cdot)$ are added on top of a linear projection. We consider three configurations. First, we set $\phi$ as a linear projection, where the dimensionality of the projected space is a hyper-parameter.[2] Then we replace batch normalization with layer normalization in the block. Finally, we add a residual link from the block's input to its output. Based on classification and regression performance across 45 datasets, we present the average performance rank of the four variants in Table 3. To avoid limiting model capacity, hyper-parameters such as the number of layers are determined based on the validation set. Further comparisons of fixed architecture configurations are listed in the appendix.

We first compare NCA with MLP vs. with the linear counterpart in Table 3. In classification tasks, MLP achieves a lower rank, highlighting the importance of incorporating nonlinearity into the model. However, in regression tasks, the linear version performs well, with MLP showing only small improvements. Although the linear projection is part of MLP's search space, the linear version benefits from a smaller hyper-parameter space, potentially resulting in better generalization.

As described in subsection 4.2, MLP uses batch normalization instead of layer normalization. Empirically, batch normalization performs better on average in both classification and regression tasks as shown in Table 3. Additionally, we compare the MLP implementation with and without residual connections. While performing similarly in classification, MLP shows superiority, especially in regression. Therefore, we adopt the MLP implementation in Table 3 for MODERNNCA.

---

[2]This "linear" version also includes the SNS sampling strategy and the nonlinear PLR encoding.

Table 4: Comparison among MODERNNCA, MLP (Gorishniy et al., 2021), and TabR (Gorishniy et al., 2024) with or without PLR encoding for numerical features. We show the change in average performance rank across the four configurations on the 45 datasets in the tiny-benchmark.

| | w/o PLR | | | w/ PLR | | |
|---|---|---|---|---|---|---|
| | MLP | TabR | MODERNNCA | MLP | TabR | MODERNNCA |
| Classification | 4.556 | 3.148 | 3.037 | 4.480 | 3.037 | 2.630 |
| Regression | 4.444 | 3.167 | 3.389 | 3.333 | 3.444 | 3.222 |



(a) classification      (b) regression

Figure 3: The change of average performance rank with different sampling rates among {10%, 30%, 50%, 80%, 100%} in SNS strategy. The dotted line denotes the rank of MODERNNCA.

**Influence of the PLR Encoding**. PLR encoding transforms numerical features into high-dimensional vectors, enhancing both model capacity and nonlinearity. To assess the impact of PLR encoding, we compare MODERNNCA with MLP and TabR, both with and without PLR encoding. Following a similar setup as in Table 3, we present the change in average performance rank across six methods in both classification and regression tasks in Table 4.

Without PLR encoding, TabR outperforms MLP, and MODERNNCA shows stronger performance in classification while performing slightly worse in regression (although still better than MLP). PLR encoding improves all methods, as evidenced by the decrease in average performance rank. In the right section of Table 4, we observe that MODERNNCA performs best in both classification and regression tasks, effectively leveraging PLR encoding better than TabR. This may be because the nonlinearity introduced by PLR compensates for the relative simplicity of MODERNNCA. The results also validate that the strength of MODERNNCA comes from a combination of its objective, architecture, and training strategy, rather than relying solely on the PLR encoding strategy.

**The Influence of Sampling Ratios**. Due to the huge computational cost of calculating distances in the learned embedding space, MODERNNCA employs a Stochastic Neighborhood Sampling (SNS) strategy, where only a subset of the training data is randomly sampled for each mini-batch. Therefore, the training time and memory cost is significantly reduced. We experiment with varying the proportion of sampled training data while keeping other hyper-parameters constant, then evaluate the corresponding test performance. As shown in Figure 3, sampling 30%-50% of the training set yields better results for MODERNNCA than using the full set. SNS not only improves training efficiency but also enhances the model's generalization ability. The plots also indicate that, with a tuned sampling ratio, MODERNNCA achieves a superior performance rank (dotted lines in the figure).

## 7 CONCLUSION

Leveraging neighborhood relationships for predictions is a classical approach in machine learning. In this paper, we revisit and enhance one of the most representative neighborhood-based methods, NCA, by incorporating modern deep learning techniques. The improved MODERNNCA establishes itself as a strong baseline for deep tabular prediction tasks, offering competitive performance while reducing the training time required to access the entire dataset. Extensive results demonstrate that MODERNNCA frequently outperforms both tree-based and deep tabular models in classification and regression tasks. Our detailed analyses shed light on the key factors driving these improvements, including the enhancements introduced to the original NCA.

# 8 REPRODUCIBILITY STATEMENT

MODERNNCA is easy to implement. The code for MODERNNCA, along with all comparison methods and datasets, is available at `https://anonymous.4open.science/r/modernNCA/`.

## REFERENCES

Nesreen K Ahmed, Amir F Atiya, Neamat El Gayar, and Hisham El-Shishiny. An empirical comparison of machine learning models for time series forecasting. *Econometric reviews*, 29(5-6): 594–621, 2010.

Takuya Akiba, Shotaro Sano, Toshihiko Yanase, Takeru Ohta, and Masanori Koyama. Optuna: A next-generation hyperparameter optimization framework. In *KDD*, 2019.

Sercan Ö. Arik and Tomas Pfister. Tabnet: Attentive interpretable tabular learning. In *AAAI*, 2021.

Lei Jimmy Ba, Jamie Ryan Kiros, and Geoffrey E. Hinton. Layer normalization. *CoRR*, abs/1607.06450, 2016.

Dara Bahri, Heinrich Jiang, Yi Tay, and Donald Metzler. Scarf: Self-supervised contrastive learning using random feature corruption. In *ICLR*, 2022.

Aurélien Bellet, Amaury Habrard, and Marc Sebban. *Metric Learning*. Morgan & Claypool Publishers, 2015.

Christopher Bishop. *Pattern recognition and machine learning*. Springer, 2006.

Vadim Borisov, Tobias Leemann, Kathrin Seßler, Johannes Haug, Martin Pawelczyk, and Gjergji Kasneci. Deep neural networks and tabular data: A survey. *IEEE Transactions on Neural Networks and Learning Systems*, abs/2110.01889:1–21, 2022.

Chun-Hao Chang, Rich Caruana, and Anna Goldenberg. NODE-GAM: neural generalized additive model for interpretable deep learning. In *ICLR*, 2022.

Jintai Chen, Kuanlun Liao, Yao Wan, Danny Z. Chen, and Jian Wu. Danets: Deep abstract networks for tabular data classification and regression. In *AAAI*, 2022.

Jintai Chen, KuanLun Liao, Yanwen Fang, Danny Chen, and Jian Wu. Tabcaps: A capsule neural network for tabular data classification with bow routing. In *ICLR*, 2023.

Jintai Chen, Jiahuan Yan, Qiyuan Chen, Danny Ziyi Chen, Jian Wu, and Jimeng Sun. Can a deep learning model be a sure bet for tabular prediction? In *KDD*, pp. 288–296, 2024.

Tianqi Chen and Carlos Guestrin. Xgboost: A scalable tree boosting system. In *KDD*, 2016.

Heng-Tze Cheng, Levent Koc, Jeremiah Harmsen, Tal Shaked, Tushar Chandra, Hrishi Aradhye, Glen Anderson, Greg Corrado, Wei Chai, Mustafa Ispir, Rohan Anil, Zakaria Haque, Lichan Hong, Vihan Jain, Xiaobing Liu, and Hemal Shah. Wide & deep learning for recommender systems. In *DLRS*, 2016.

Jason V. Davis, Brian Kulis, Prateek Jain, Suvrit Sra, and Inderjit S. Dhillon. Information-theoretic metric learning. In *ICML*, 2007.

Manuel Fernández Delgado, Eva Cernadas, Senén Barro, and Dinani Gomes Amorim. Do we need hundreds of classifiers to solve real world classification problems? *Journal of Machine Learning Research*, 15(1):3133–3181, 2014.

Janez Demsar. Statistical comparisons of classifiers over multiple data sets. *Journal of Machine Learning Research*, 7:1–30, 2006.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: pre-training of deep bidirectional transformers for language understanding. In *NAACL-HLT*, 2019.

Nicholas Frosst, Nicolas Papernot, and Geoffrey Hinton. Analyzing and improving representations with the soft nearest neighbor loss. In *International conference on machine learning*, pp. 2012–2020. PMLR, 2019a.

Nicholas Frosst, Nicolas Papernot, and Geoffrey E. Hinton. Analyzing and improving representations with the soft nearest neighbor loss. In *ICML*, volume 97, pp. 2012–2020, 2019b.

Amir Globerson and Sam T. Roweis. Metric learning by collapsing classes. In *NIPS*, pp. 451–458, 2005.

Jacob Goldberger, Sam T. Roweis, Geoffrey E. Hinton, and Ruslan Salakhutdinov. Neighbourhood components analysis. In *NIPS*, 2004.

Yury Gorishniy, Ivan Rubachev, Valentin Khrulkov, and Artem Babenko. Revisiting deep learning models for tabular data. In *NeurIPS*, 2021.

Yury Gorishniy, Ivan Rubachev, and Artem Babenko. On embeddings for numerical features in tabular deep learning. In *NeurIPS*, 2022.

Yury Gorishniy, Ivan Rubachev, Nikolay Kartashev, Daniil Shlenskii, Akim Kotelnikov, and Artem Babenko. Tabr: Tabular deep learning meets nearest neighbors in 2023. In *ICLR*, 2024.

Léo Grinsztajn, Edouard Oyallon, and Gaël Varoquaux. Why do tree-based models still outperform deep learning on typical tabular data? In *NeurIPS*, 2022.

Huifeng Guo, Ruiming Tang, Yunming Ye, Zhenguo Li, and Xiuqiang He. Deepfm: A factorization-machine based neural network for CTR prediction. In *IJCAI*, 2017.

Md. Rafiul Hassan, Sadiq Al-Insaif, Muhammad Imtiaz Hossain, and Joarder Kamruzzaman. A machine learning approach for prediction of pregnancy outcome following IVF treatment. *Neural Computing and Applications*, 32(7):2283–2297, 2020.

Noah Hollmann, Samuel Müller, Katharina Eggensperger, and Frank Hutter. Tabpfn: A transformer that solves small tabular classification problems in a second. In *ICLR*, 2023.

Cheng-Kang Hsieh, Longqi Yang, Yin Cui, Tsung-Yi Lin, Serge J. Belongie, and Deborah Estrin. Collaborative metric learning. In *WWW*, 2017.

Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *ICML*, 2015.

Alan Jeffares, Tennison Liu, Jonathan Crabbé, Fergus Imrie, and Mihaela van der Schaar. Tangos: Regularizing tabular neural networks through gradient orthogonalization and specialization. In *ICLR*, 2023.

Arlind Kadra, Marius Lindauer, Frank Hutter, and Josif Grabocka. Well-tuned simple nets excel on tabular datasets. In *NeurIPS*, pp. 23928–23941, 2021.

Liran Katzir, Gal Elidan, and Ran El-Yaniv. Net-dnf: Effective deep modeling of tabular data. In *ICLR*, 2021.

Guolin Ke, Qi Meng, Thomas Finley, Taifeng Wang, Wei Chen, Weidong Ma, Qiwei Ye, and Tie-Yan Liu. Lightgbm: A highly efficient gradient boosting decision tree. In *NIPS*, 2017.

Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron Maschinot, Ce Liu, and Dilip Krishnan. Supervised contrastive learning. In *NeurIPS*, 2020.

Brian Kulis. Metric learning: A survey. *Foundations and Trends in Machine Learning*, 5(4), 2013.

Steinar Laenen and Luca Bertinetto. On episodes, prototypical networks, and few-shot learning. In *NeurIPS*, pp. 24581–24592, 2021.

Dong C Liu and Jorge Nocedal. On the limited memory bfgs method for large scale optimization. *Mathematical programming*, 45(1):503–528, 1989.

Duncan C. McElfresh, Sujay Khandagale, Jonathan Valverde, Vishak Prasad C., Ganesh Ramakrishnan, Micah Goldblum, and Colin White. When do neural nets outperform boosted trees on tabular data? In *NeurIPS*, 2023.

Martin Renqiang Min, Laurens van der Maaten, Zineng Yuan, Anthony J. Bonner, and Zhaolei Zhang. Deep supervised t-distributed embedding. In *ICML*, pp. 791–798, 2010.

Mehryar Mohri, Afshin Rostamizadeh, and Ameet Talwalkar. *Foundations of Machine Learning*. MIT Press, 2012.

Lennart J Nederstigt, Steven S Aanen, Damir Vandic, and Flavius Frasincar. Floppies: a framework for large-scale ontology population of product information from tabular data in e-commerce stores. *Decision Support Systems*, 59:296–311, 2014.

Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32, 2019.

F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.

Sergei Popov, Stanislav Morozov, and Artem Babenko. Neural oblivious decision ensembles for deep learning on tabular data. In *ICLR*, 2020.

Liudmila Ostroumova Prokhorenkova, Gleb Gusev, Aleksandr Vorobev, Anna Veronika Dorogush, and Andrey Gulin. Catboost: unbiased boosting with categorical features. In *NeurIPS*, 2018.

Matthew Richardson, Ewa Dominowska, and Robert Ragno. Predicting clicks: estimating the click-through rate for new ads. In *WWW*, 2007.

Ivan Rubachev, Artem Alekberov, Yury Gorishniy, and Artem Babenko. Revisiting pretraining objectives for tabular deep learning. *CoRR*, abs/2207.03208, 2022.

Ruslan Salakhutdinov and Geoffrey E. Hinton. Learning a nonlinear embedding by preserving class neighbourhood structure. In *AISTATS*, volume 2, pp. 412–419, 2007.

Florian Schroff, Dmitry Kalenichenko, and James Philbin. Facenet: A unified embedding for face recognition and clustering. In *CVPR*, 2015.

Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *ICLR*, 2015.

Kihyuk Sohn. Improved deep metric learning with multi-class n-pair loss objective. In *NIPS*, 2016.

Gowthami Somepalli, Avi Schwarzschild, Micah Goldblum, C. Bayan Bruss, and Tom Goldstein. SAINT: Improved neural networks for tabular data via row attention and contrastive pre-training. In *NeurIPS Workshop*, 2022.

Hyun Oh Song, Yu Xiang, Stefanie Jegelka, and Silvio Savarese. Deep metric learning via lifted structured feature embedding. In *CVPR*, 2016.

Nitish Srivastava, Geoffrey E. Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: a simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research*, 15(1):1929–1958, 2014.

Daniel Tarlow, Kevin Swersky, Laurent Charlin, Ilya Sutskever, and Richard S. Zemel. Stochastic k-neighborhood selection for supervised and unsupervised learning. In *ICML*, pp. 199–207, 2013.

Talip Ucar, Ehsan Hajiramezanali, and Lindsay Edwards. Subtab: Subsetting features of tabular data for self-supervised representation learning. In *NeurIPS*, pp. 18853–18865, 2021.

Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of machine learning research*, 9(11), 2008.

Joaquin Vanschoren, Jan N Van Rijn, Bernd Bischl, and Luis Torgo. Openml: networked science in machine learning. *ACM SIGKDD Explorations Newsletter*, 15(2):49–60, 2014.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NIPS*, 2017.

Jarkko Venna, Jaakko Peltonen, Kristian Nybo, Helena Aidos, and Samuel Kaski. Information retrieval perspective to nonlinear dimensionality reduction for data visualization. *Journal of Machine Learning Research*, 11:451–490, 2010.

Oriol Vinyals, Charles Blundell, Tim Lillicrap, Koray Kavukcuoglu, and Daan Wierstra. Matching networks for one shot learning. In *NIPS*, 2016.

Ruoxi Wang, Rakesh Shivanna, Derek Zhiyuan Cheng, Sagar Jain, Dong Lin, Lichan Hong, and Ed H. Chi. DCN V2: improved deep & cross network and practical lessons for web-scale learning to rank systems. In *WWW*, 2021.

Tianjun Wei, Jianghong Ma, and Tommy W. S. Chow. Collaborative residual metric learning. In *SIGIR*, 2023.

Kilian Q. Weinberger and Lawrence K. Saul. Distance metric learning for large margin nearest neighbor classification. *Journal of Machine Learning Research*, 10:207–244, 2009.

Jing Wu, Suiyao Chen, Qi Zhao, Renat Sergazinov, Chen Li, Shengjie Liu, Chongchao Zhao, Tianpei Xie, Hanqing Guo, Cheng Ji, Daniel Cociorva, and Hakan Brunzell. Switchtab: Switched autoencoders are effective tabular learners. In *AAAI*, pp. 15924–15933, 2024.

Eric P. Xing, Andrew Y. Ng, Michael I. Jordan, and Stuart Russell. Distance metric learning with application to clustering with side-information. In *NIPS*, 2002.

Xun Yang, Meng Wang, and Dacheng Tao. Person re-identification with metric learning using privileged information. *IEEE Transactions on Image Processing*, 27(2):791–805, 2018.

Han-Jia Ye, Si-Yang Liu, Hao-Run Cai, Qi-Le Zhou, and De-Chuan Zhan. A closer look at deep learning on tabular data. *CoRR*, abs/2407.00956, 2024a.

Hangting Ye, Wei Fan, Xiaozhuang Song, Shun Zheng, He Zhao, Dan dan Guo, and Yi Chang. Ptarl: Prototype-based tabular representation learning via space calibration. In *ICLR*, 2024b.

Dong Yi, Zhen Lei, Shengcai Liao, and Stan Z. Li. Deep metric learning for person re-identification. In *ICME*, 2014.

The Appendix consists of two sections:

- Appendix A: Datasets and implementation details.
- Appendix B: Additional experimental results.

## APPENDIX A  DATASETS AND IMPLEMENTATION DETAILS

In this section, we outline the preprocessing steps applied to the datasets before training, as well as descriptions of the datasets used.

### A.1  DATA PRE-PROCESSING

We follow the data preprocessing pipeline from Gorishniy et al. (2021) for all methods. For numerical features, we apply standardization by subtracting the mean and scaling the values. For categorical features, we use one-hot encoding to convert them for model input.

### A.2  DATASET INFORMATION

We use the recent large-scale tabular benchmark from Ye et al. (2024a), which includes 300 datasets covering various domains such as healthcare, biology, finance, education, and physics. The dataset sizes range from 1,000 to 1 million instances. More detailed information on the datasets can be found in Ye et al. (2024a).

For each dataset, we randomly sample 20% of the instances to form the test set. The remaining 80% is split further, with 20% of which held out as a validation set. The validation set is used to tune hyper-parameters and apply early stopping. The hyper-parameters with which the model performs best on the validation set are selected for final evaluation with the test set.

The datasets used in our analyses and ablation studies follow the tiny-benchmark in Ye et al. (2024a), which consists of 45 datasets. The performance rankings of methods on this smaller benchmark are consistent with those on the full benchmark, making it a useful probe for tabular analysis.

### A.3  HARDWARE

The majority of experiments, including those measuring time and memory overhead, were conducted on a Tesla V100 GPU.

### A.4  POTENTIAL ALTERNATIVE IMPLEMENTATION

We explore an alternative strategy to learn the embedding $\phi$ in two steps. First, we apply Supervised Contrastive loss (Sohn, 2016; Khosla et al., 2020), where supervision is generated within a mini-batch. After learning $\phi$, we use KNN for classification or regression during inference. In the regression scenario, label values are discretized, and we refer to this baseline method as Tabular Contrastive (TabCon). Empirically, we find that certain components of MODERNNCA, such as the Soft-NN loss for prediction, cannot be directly applied to TabCon, even when $\phi$ is implemented using the same nonlinear MLP as in MODERNNCA. Despite this, the TabCon baseline remains competitive with FT-Transformer (FT-T), achieving average ranks similar to L-NCA in both classification and regression tasks.

## APPENDIX B  ADDITIONAL EXPERIMENTS

### B.1  VISUALIZATION RESULTS

To better analyze the properties of MODERNNCA, we visualize the learned embeddings $\phi(\boldsymbol{x})$ of MODERNNCA, TabCon (mentioned in subsection A.4), and TabR using TSNE (Van der Maaten & Hinton, 2008). As shown in Figure 4, all deep tabular methods transform the embedding spaces to be more helpful for classification or regression compared to the raw features. The embedding space learned by TabCon clusters samples of the same class together and separates samples of different

(a) AD ↑ Raw Feature    (b) AD ↑ TabR    (c) AD ↑ MODERNNCA    (c) AD ↑ TabCon

(a) PH ↑ Raw Feature    (b) PH ↑ TabR    (c) PH ↑ MODERNNCA    (c) PH ↑ TabCon

(a) CA↓ Raw Feature    (b) CA↓ TabR    (c) CA↓ MODERNNCA    (c) CA↓ TabCon

(a) MIA↓ Raw Feature    (b) MIA↓ TabR    (c) MIA↓ MODERNNCA    (c) MIA↓ TabCon

Figure 4: Visualization of the embedding space of different methods.

classes, often clustering same-class instances into a single cluster. However, it still struggles with some hard-to-distinguish samples. TabR and MODERNNCA, on the other hand, divide samples of the same class into multiple clusters, ensuring that similar samples are positioned closer to each other. This strategy aligns with the prediction mechanism of KNN, where good performance is achieved by clustering instances with similar neighbors together rather than into a single cluster. The embedding space learned by MODERNNCA is more discriminative than that learned by TabR. The main reason is that TabR leverages an additional architecture to modify the prediction score for each instance, making the learned embedding space less discriminative compared to MODERNNCA.

### B.2 ADDITIONAL ABLATION STUDIES

**The Influence of Distance Functions**. The predicted label of a target instance $\boldsymbol{x}_i$ is determined by the label of its neighbors in the learned embedding space projected by $\phi$. The distance function $\mathrm{dist}(\cdot, \cdot)$ is the key to determining the pairwise relationship between instances in the embedding space and influences the weights in Equation 5.

In MODERNNCA, we choose Euclidean distance

$$\mathrm{dist}_{\mathrm{EUC}}(\phi(\boldsymbol{x}_i), \phi(\boldsymbol{x}_j)) = \sqrt{(\phi(\boldsymbol{x}_i) - \phi(\boldsymbol{x}_j))^\top (\phi(\boldsymbol{x}_i) - \phi(\boldsymbol{x}_j))} = \|\phi(\boldsymbol{x}_i) - \phi(\boldsymbol{x}_j)\|_2 . \quad (7)$$

We also utilize other distance functions, *e.g.*, the squared Euclidean distance, $\mathrm{dist}_{\mathrm{EUC}}^2(\phi(\boldsymbol{x}_i), \phi(\boldsymbol{x}_j))$, the $\ell_1$-norm distance

$$\mathrm{dist}(\phi(\boldsymbol{x}_i), \phi(\boldsymbol{x}_j)) = \|\phi(\boldsymbol{x}_i) - \phi(\boldsymbol{x}_j)\|_1 , \quad (8)$$

the (negative) cosine similarity $\mathrm{dist}(\phi(\boldsymbol{x}_i), \phi(\boldsymbol{x}_j)) = -(\boldsymbol{x}_i^\top \boldsymbol{x}_j)/(\|\boldsymbol{x}_i\|_2 \|\boldsymbol{x}_j\|_2)$, and the (negative) inner product $\mathrm{dist}(\phi(\boldsymbol{x}_i), \phi(\boldsymbol{x}_j)) = -\phi(\boldsymbol{x}_i)^\top \phi(\boldsymbol{x}_j)$. The results using different distance functions

Table 5: Comparison among various distances used to implement Equation 5, where Euclid distance is the default choice in MODERNNCA. We show the change in average performance rank (lower is better) across the five configurations on the 45 datasets in the tiny-benchmark.

|  | Euclid | Dot Product | Cosine | Squared Euclid | L1-Norm |
|---|---|---|---|---|---|
| Classification | 2.593 | 3.852 | 2.111 | 2.630 | 3.769 |
| Regression | 2.500 | 3.222 | 2.529 | 2.722 | 3.889 |

Table 6: Comparison of different loss functions. The log loss used in MODERNNCA, the original NCA's summation loss, the MCML loss, and the t-distribution loss. The change in average performance rank (lower is better) is presented across these four configurations on the 45 datasets in the tiny-benchmark.

|  | MODERNNCA | NCA | MCML | t-distribution |
|---|---|---|---|---|
| Classification | 2.074 | 2.519 | 3.074 | 2.333 |
| Regression | 1.500 | - | - | 1.540 |

are listed in Table 5, which contains the average performance rank over 45 datasets among the five variants. On average, Euclidean distance performs well across both classification and regression tasks. While cosine distance yields better results on classification datasets (with an average performance rank of 4.5939 compared to MODERNNCA and 20 other methods across 300 datasets, please check Figure 2 for details), its advantage diminishes on regression tasks.

**Other Possible Loss Functions**. NCA (Goldberger et al., 2004) originally explored two loss functions: one that maximizes the sum of probabilities in Equation 3, and another that minimizes the sum of log probabilities as in Equation 1. The former was selected in the original implementation of NCA due to its better performance. We also investigated several alternative loss functions for NCA. For instance, MCML (Globerson & Roweis, 2005) minimizes the KL-divergence between the learned embedding in Equation 2 and a constructed ground-truth label distribution for each instance, but it only applies to classification tasks. Another variant is the t-distributed NCA (Min et al., 2010), which uses a heavy-tailed t-distribution to measure pairwise similarities in the objective function. We tested both MCML and the t-distribution loss functions in MODERNNCA, and the results are summarized in Table 6, showing the average ranks across 45 datasets. The log objective in Equation 1 performs best for classification tasks and slightly outperforms the t-distribution variant in regression tasks.

**The Influence of Sampling Strategy**. As mentioned before, SNS randomly samples a subset of training data for each mini-batch when calculating the loss of Equation 5. We also investigate whether we could further improve the classification/regression ability of the model when we incorporate richer information during the sampling process, *e.g.*, the label of the instances.

We consider two other sampling strategies in addition to the fully random one we used before. First is class-wise random sampling, which means that given a proportion, we sample from each class in the training set and combine them together. This strategy takes advantage of the training label information and keeps the instances from all classes that will exist in the sampled subset. Besides, we also consider the sampling strategy based on the pairwise distances between instances. Since the neighbors of an instance may contribute more (with larger weights) in Equation 5, so given a mini-batch, we first calculate the Euclidean distance between instances in the batch and all the training set with the embedding function $\phi$ in the current epoch. Then we sample the training set based on the reciprocal of the pairwise distance value. In detail, given an instance $x_i$, we provide instance-specific neighborhood candidates and $x_j$ in the training set is sampled based on the probability $\sim 1/(\text{dist}(\phi(x_i), \phi(x_j)))^\tau$. $\tau$ is a non-negative hyper-parameter to calibrate the distribution. The distance calculation requires forward passes of the model $\phi$ over all the training instances, and the instance-specific neighborhood makes the loss related to a wide range of the training data. Therefore, the distance-based sampling strategy has a low training speed and high computational burden.

The comparison results, *i.e.*, the average performance rank, among different sampling strategies on 45 datasets are listed in Table 7. We empirically find the label-based sampling strategy cannot provide

17

Table 7: Comparison of different sampling strategies: "Random", "Label", and "Distance" represent MODERNNCA's naive uniform sampling, class-wise random sampling, and distance-based sampling, respectively. The change in average performance rank (lower is better) is presented across these three configurations on the 45 datasets in the tiny-benchmark.

|  | Random | Label | Distance |
|---|---|---|---|
| Classification | 1.869 | 2.230 | 1.901 |
| Regression | 1.508 | - | 1.492 |

Table 8: Comparison of various architecture choices based on a fixed 2-layer MLP. We only tune architecture-independent hyper-parameters for different variants. The change in average performance rank (lower is better) is shown across three configurations (default, Layer Norm, and Residual) on the 45 datasets in the tiny-benchmark.

|  | MLP | w/ LayerNorm | ResNet |
|---|---|---|---|
| Classification | 1.905 | 2.048 | 2.048 |
| Regression | 1.813 | 2.313 | 1.875 |

further improvements. Although the distance-based strategy helps in certain cases, the improvements are limited. Taking a holistic consideration of the performance and efficiency, we choose to use the vanilla random sampling in MODERNNCA.

**Comparison between Different Deep Architectures**. Unlike the ablation studies in subsection 6.2, where we fixed the model family and tuned detailed hyper-parameters (such as the number of layers and network width) based on the validation set, here we fix the main architecture as a two-layer MLP and only tune architecture-independent hyper-parameters, such as the learning rate.

With this base MLP architecture, we evaluate three variants: the base MLP, one with batch normalization replaced by layer normalization, and one with an added residual link. The average ranks of the three variants across 45 datasets are presented in Table 8. We observe that the basic MLP remains a better choice compared to the versions with a residual link or layer normalization.

## B.3 RUN-TIME AND MEMORY USAGE ESTIMATION

We make a run-time and memory usage comparison in Figure 1. Here are the steps that we take to perform the estimation. First, we tuned all models on the validation set for 100 iterations, saving the optimal parameters ever found. Next, we ran the models for 15 iterations with the tuned parameters and saved the best checkpoint on the validation set. The run-time for the models was estimated using the average time taken by the tuned model to run one seed in the training and validation stage.

We present the average results of run-time and memory usage estimation across the full benchmark (300 datasets) in Table 9.

Table 9: Training time and memory usage estimation for different tuned models over 300 datasets. The average rank represents the mean performance ranking of these models based on the performance metrics (RMSE for regression and accuracy for classification).

| Model | M-NCA | L-NCA | MLP | MLP-PLR | FT-T | TabR | XGBoost | CatBoost |
|---|---|---|---|---|---|---|---|---|
| Training Time (s) | 87.5 | 33.62 | 30.36 | 42.87 | 111.91 | 173.34 | 4.53 | 20.48 |
| Memory Usage (GB) | 5.36 | 1.42 | 1.15 | 2.37 | 4.98 | 10.13 | 0.84 | 1.06 |
| Average Rank | 4.56 | 6.30 | 7.53 | 6.94 | 6.29 | 5.36 | 5.62 | 4.61 |

## B.4 FULL RESULT ON THE BENCHMARK

Table 10: Detailed results (average accuracy) for all classification datasets over all methods. Due to space constraints, dataset IDs from Ye et al. (2024a) are used in place of dataset names. "XGB", "RF", "LG", "RN", "ND", "ST", "TG", "DAN", "FTT", "DCN", "TT", "PT", "EF", and "TR" denote "XGBoost", "Random Forest", "LightGBM", "ResNet", "NODE", "SwitchTab", "TANGOS", "DANets", "FT-T", "DCNv2", "PTaRL", "Excelformer" and "TabR", respectively.

| ID | M-NCA | L-NCA | TR | XGB | CGB | MLP | FTT | DCN | ND | MP | EF | DAN | TG | TC | ST | LG | RF | SVM | KNN | SW | PT |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 5 | .3871 | .3907 | .3874 | .4301 | .3815 | .3902 | .3678 | .3588 | .2780 | .3807 | .2901 | .3699 | .3480 | .3701 | - | .4368 | .4108 | .3465 | .3707 | .3514 | .3596 |
| 8 | .9437 | .9435 | .9426 | .9428 | .9467 | .9432 | .9420 | .9421 | .9422 | - | - | .9425 | .9421 | .9422 | .9439 | .9447 | .9426 | .9422 | .9411 | .9421 | .9431 |
| 10 | .7545 | .7453 | .7569 | .7455 | .7445 | .7344 | .7339 | .7343 | .7297 | .7337 | .7311 | .7329 | .7358 | .7399 | .7422 | .7366 | .7596 | .7244 | .7321 | | |
| 11 | .9859 | .9833 | .9868 | .9876 | .9878 | .9846 | .9860 | .9841 | .9856 | .9852 | .9856 | .9840 | .9838 | .9834 | - | .9873 | .9852 | .9781 | .9837 | .9818 | .9840 |
| 12 | .5832 | .5770 | .5792 | .5838 | .5893 | .5780 | .5850 | .5838 | .5881 | .5786 | .5870 | .5779 | .5819 | .5810 | - | .5837 | .5774 | .5185 | .5629 | .5624 | .5475 |
| 17 | .8108 | .8035 | .8045 | .8099 | .8133 | .8003 | .8131 | .8153 | .8117 | - | .8134 | .8009 | .7992 | .7981 | - | .8075 | .7979 | .6407 | .7779 | .7703 | .8112 |
| 18 | .8567 | .8567 | .8762 | .8688 | .8759 | .8626 | .8709 | .8674 | .8654 | .8630 | .8706 | .8636 | .8659 | .8650 | .8677 | .8676 | .8681 | .8073 | .8485 | .8495 | .8646 |
| 19 | .6973 | .6749 | .7050 | .6858 | .6761 | .6975 | .7015 | .7012 | .6542 | .6930 | .6970 | .7075 | .7005 | .7052 | .6995 | .6853 | .6891 | .7082 | .6679 | .6903 | .7085 |
| 23 | .8676 | .8665 | .8674 | .8676 | .8678 | .8680 | .8676 | .8681 | .8678 | .8680 | .8678 | .8680 | .8677 | .8679 | .8671 | .8678 | .8665 | .8642 | .8641 | .8674 | .8681 |
| 27 | .9053 | .8961 | .9135 | .9000 | .9082 | .8847 | .8923 | .8858 | .8730 | .8878 | .8915 | .8830 | .8835 | .8812 | .8853 | .9070 | .8765 | .8413 | .8522 | .8135 | .8805 |
| 28 | .7320 | .7298 | .7314 | .7331 | .7340 | .7339 | .7322 | .7321 | .7320 | .7326 | .7317 | .7323 | .7298 | .7316 | - | .7335 | .7325 | .6893 | .6775 | .7284 | .7321 |
| 29 | .8331 | - | .8322 | .8325 | .8326 | .8334 | .8333 | .8327 | .8331 | .8333 | .8336 | .8321 | .8323 | .8333 | .8334 | .8332 | .8325 | .7921 | .8325 | .8325 | .8331 |
| 36 | .7911 | .8003 | .8202 | .7814 | .7972 | .7315 | .7770 | .7771 | .7194 | .7697 | - | .7101 | .7230 | .7130 | .7555 | .7751 | .7269 | .6311 | .7786 | .6902 | .7754 |
| 37 | .8914 | .8832 | .8609 | .8754 | .8905 | .8787 | .8927 | .8917 | .8580 | .8935 | .8850 | .8679 | .8743 | .8854 | .8932 | .8783 | .8801 | .8815 | .8705 | .8662 | .8876 |
| 39 | .9677 | - | - | .9689 | .9681 | .9648 | .9663 | .9634 | .9582 | .9665 | - | .9523 | .9624 | .9569 | .9639 | .9666 | .9529 | .9370 | .9161 | .9602 | |
| 40 | .7469 | .6987 | .7576 | .7261 | .7414 | .7391 | .7319 | .7671 | .5772 | .7460 | .7380 | .7388 | .7391 | .7065 | .7348 | .7328 | .7351 | .7227 | .6494 | .6349 | .7322 |
| 42 | .6972 | .6937 | .6710 | .6723 | .6672 | .6692 | .6728 | .6627 | .6634 | .6692 | .6694 | .6619 | .6640 | .6666 | .6675 | .6779 | .6650 | .6286 | .6455 | .6597 | .6655 |
| 43 | .8459 | .8401 | .8384 | .8462 | .8500 | .8458 | .8477 | .8497 | .8438 | .8470 | .8596 | .8489 | .8510 | .8412 | .8505 | .8362 | .8517 | .7079 | .8195 | .7922 | .8379 |
| 44 | .7428 | .7340 | .7396 | .7411 | .7427 | .7391 | .7410 | .7374 | .7451 | .7445 | .7411 | .7361 | .7356 | .7411 | .7371 | .7424 | .7390 | .7443 | .7251 | .7319 | .7340 |
| 45 | .7057 | .7301 | .7428 | .6714 | .6816 | .7384 | .7346 | .7370 | .5577 | .6945 | .6807 | .7277 | .7315 | .7322 | .7012 | .6812 | .6750 | .7391 | .6213 | .6625 | .7341 |
| 46 | .7115 | .7116 | .6381 | .6630 | .7121 | .6901 | .7108 | .5900 | .7008 | .7052 | .7112 | .6725 | .7053 | - | .6642 | .5849 | .7034 | .5765 | .6539 | .7123 | |
| 47 | .7544 | .7290 | .7520 | .6627 | .6398 | .6901 | .5179 | .7480 | .6523 | .6500 | .5297 | .5370 | .5010 | .6111 | .5177 | .6618 | .6554 | .7520 | .5940 | .5232 | .6191 |
| 48 | .7526 | .7571 | .7655 | .6487 | .6558 | .7304 | .6703 | .7450 | .5506 | .6774 | .6421 | .6716 | .7052 | .6957 | .6153 | .6460 | .6416 | .7664 | .5913 | .6118 | .7199 |
| 49 | .7209 | .7214 | .7202 | .6407 | .6495 | .7173 | .6913 | .7132 | .5436 | .6777 | .6428 | .7058 | .7129 | .6875 | - | .6440 | .5818 | .7142 | .5480 | .5387 | .7088 |
| 50 | .7109 | .7030 | .7080 | .6023 | .6158 | .6985 | .6184 | .6909 | .5225 | .6453 | .6146 | .6562 | .5270 | .6540 | - | .6228 | .5653 | .7041 | .5254 | .5146 | .6512 |
| 51 | .7098 | .7072 | .7109 | .6080 | .6165 | .6831 | .6717 | .6976 | .5347 | .6162 | .5690 | .6494 | .5362 | .6216 | - | .6225 | .5657 | .6913 | .5313 | .5329 | .6923 |
| 52 | .7041 | .7134 | .7124 | .6025 | .5889 | .7204 | .6665 | .7152 | .7063 | .6136 | .6265 | .7014 | .6966 | .6946 | .6639 | .5913 | .5713 | .7139 | .5204 | .5537 | .5955 |
| 53 | .7118 | .6895 | .6432 | .6121 | .6259 | .7061 | .6880 | .7045 | .5338 | .6379 | .6644 | .6871 | .5885 | .6838 | - | .6272 | .5718 | .7099 | .5461 | .5334 | .7046 |
| 55 | .8719 | .8519 | .8690 | .8569 | .8597 | .8650 | .8711 | .8639 | .8142 | - | - | .8587 | .8646 | .8661 | .8592 | .8561 | .8077 | .7765 | .8046 | .8294 | .8615 |
| 56 | .7767 | .7753 | .7871 | .7500 | .7811 | .7822 | .7702 | .7693 | .7862 | .7904 | - | .7864 | .7882 | .7924 | .7753 | .7733 | .7738 | .7800 | .7900 | .7720 | .7816 |
| 58 | .5990 | .6281 | .5894 | .6096 | .6808 | .5846 | .4931 | .5890 | .5696 | - | .5633 | .5710 | .6012 | .6010 | .6096 | .5579 | .5656 | .5000 | .5377 | .5131 | |
| 59 | .6571 | .6712 | .6685 | .6665 | .6796 | .6400 | .5325 | .6637 | .5329 | - | - | .6373 | .6279 | .6704 | .6600 | .6823 | .6502 | .4979 | .5469 | .5427 | .5981 |
| 60 | .5925 | .5724 | .5914 | .5987 | .5926 | .6009 | .5986 | .5895 | .5917 | .5948 | .5941 | .5900 | .5956 | .5982 | .5948 | .5882 | .5959 | .5888 | .5726 | .5868 | .5966 |
| 61 | .6951 | .5885 | .7392 | .6894 | .6928 | .6342 | .6565 | .6250 | .5230 | .6679 | .6844 | .6134 | .6105 | .6427 | - | .6833 | .6214 | .4558 | .6567 | .5235 | .6016 |
| 63 | .7958 | .7961 | .7906 | .7974 | .7951 | .7835 | .7980 | .7972 | .7910 | .7928 | .7940 | .7855 | .7837 | .7871 | - | .7967 | .7939 | .7616 | .7811 | .7753 | .7977 |
| 65 | .8653 | .8689 | .8766 | .8478 | .8488 | .8626 | .8723 | .8746 | .8408 | .8760 | .8760 | .8565 | .8667 | .8467 | .8537 | .8683 | .8401 | .8540 | .8449 | | |
| 67 | .8970 | .8938 | .8933 | .8972 | .8921 | .8739 | .8850 | .8746 | .8443 | .8750 | .8825 | .8729 | .8735 | .8728 | .8761 | .8962 | .8821 | .8062 | .8707 | .8365 | .8779 |
| 68 | .9685 | .9713 | .9650 | .9232 | .9207 | .9301 | .9335 | .9343 | .7264 | .9329 | .9339 | .9154 | .9207 | .9377 | - | .8427 | .8905 | .8987 | .8950 | .7736 | .9262 |
| 72 | .9956 | .9980 | .9959 | .9803 | .9830 | .9898 | .9869 | .9879 | .9481 | .9870 | .9821 | .9876 | .9870 | .9796 | .9883 | .9821 | .9605 | .9282 | .9860 | .9362 | .9866 |
| 75 | .7999 | .7617 | .7854 | .8163 | .8103 | .7663 | .8038 | .7565 | .7741 | .8099 | .8070 | .7596 | .7552 | .7676 | .7772 | .8101 | .8134 | .7388 | .7193 | .7516 | .7606 |
| 83 | .8782 | .9036 | .8929 | .9105 | .8871 | .8828 | .8968 | .8909 | .8491 | .8911 | .8772 | .8788 | .8800 | .8705 | .8479 | .9026 | .8836 | .8879 | .8576 | .8693 | .8848 |
| 84 | .8790 | .8608 | .8813 | .8763 | .8787 | .8724 | .8759 | .8718 | .8701 | .8705 | .8758 | .8691 | .8691 | .8725 | .8732 | .8788 | .8690 | .7934 | .8465 | .8512 | .8698 |
| 85 | .9037 | .8990 | .8772 | .8827 | .8805 | .8704 | .8688 | .8643 | .8615 | .8829 | .9012 | .8845 | .8760 | .8862 | .8738 | .8802 | .8763 | .8708 | .8728 | .8754 | .8808 |
| 88 | .9857 | .9747 | .9575 | .9090 | .9357 | .8672 | .9367 | .9550 | .8953 | .9637 | .9527 | .7978 | .9115 | .9250 | .9025 | .9098 | .8760 | .8348 | .7450 | .7852 | .9477 |
| 91 | .8326 | .8456 | .8344 | .8423 | .8427 | .8303 | .8363 | .8300 | .8425 | .8349 | .8330 | .8298 | .8358 | .8355 | .8364 | .8507 | .8319 | .8289 | .8266 | .8386 | |
| 95 | .9767 | .9766 | .9760 | .9731 | .9679 | .9703 | .9676 | .9696 | .9486 | - | .9697 | .9699 | .9708 | .9681 | - | .9734 | .9624 | .9290 | .9584 | .9445 | .9698 |
| 97 | .8799 | .8824 | .8744 | .8753 | .8750 | .8765 | .8710 | .8586 | .8747 | .8725 | .8744 | .8759 | .8756 | .8744 | .8759 | .8694 | .8731 | .8889 | .8750 | .8670 | .8747 |
| 98 | .7524 | .7351 | .7597 | .7645 | .7775 | .7606 | .7558 | .7364 | .7706 | .7576 | .7576 | .7251 | .7710 | .7524 | .7675 | .7814 | .7727 | .7468 | .7143 | .7398 | .7580 |
| 99 | .8804 | .8823 | .8772 | .8561 | .8775 | .8855 | .8826 | .8762 | .8801 | .8689 | .8746 | .8817 | .8813 | .8804 | .8778 | .8807 | .8743 | .9027 | .8756 | .8797 | .8797 |
| 100 | .8699 | .8675 | .8687 | .8716 | .8725 | .8691 | .8672 | .8717 | .8607 | .8701 | .8685 | .8508 | .8661 | .8663 | .8697 | .8747 | .8740 | .8587 | .8640 | .8513 | .8651 |
| 101 | .8455 | .8755 | .8603 | .8506 | .8730 | .8667 | .8695 | .8798 | .8686 | .6976 | .8773 | .8682 | .8408 | .8742 | .8442 | .8793 | .8682 | .8408 | .8742 | .8910 | .8610 | .8537 |
| 102 | .8635 | .8532 | .8749 | .8528 | .8689 | .8448 | .8566 | .8591 | .8474 | .8426 | - | .8433 | .8432 | .8432 | .8490 | .8528 | .8401 | .8221 | .8331 | .8318 | .8595 |
| 103 | .9738 | - | .9742 | .9768 | .9752 | .9728 | .9703 | .9683 | .9728 | .9732 | .9740 | .9710 | .9725 | .9726 | - | .9764 | .9756 | .9342 | .9233 | .7779 | .9664 |
| 104 | .9918 | .9903 | .9918 | .9921 | .9920 | .9913 | .9920 | .9911 | .9890 | .9920 | .9904 | .9922 | .9920 | .9915 | .9915 | .9900 | .9913 | .9912 | .9912 | .9913 | .9913 |
| 107 | .7392 | - | - | .7387 | .7399 | .7390 | .7403 | .7391 | .7401 | .7397 | .7398 | .7390 | .7382 | .7391 | .7390 | .7383 | .7308 | .7234 | .7236 | .7330 | .7387 |
| 111 | .9456 | .9370 | .9418 | .9557 | .9517 | .9162 | .9519 | .9196 | .9129 | .9250 | .9320 | .9231 | .9225 | .9172 | .9337 | .9539 | .9544 | .8756 | .8951 | .8865 | .9173 |
| 112 | .6203 | .5814 | .4549 | .3177 | .3531 | .4659 | .4164 | .4755 | .1709 | .4379 | - | .4530 | .3244 | .4215 | .3399 | .3458 | .1795 | .4103 | .3055 | .3438 | |
| 115 | .6280 | .6235 | .6263 | .6144 | .6250 | .6301 | .6310 | .6258 | .6505 | .6233 | - | .6258 | .6271 | .6309 | .6277 | .6286 | .6293 | .6098 | .6387 | .6143 | .6313 |
| 117 | .9843 | .9817 | .9872 | .9883 | .9879 | .9885 | .9849 | .9851 | .9461 | .9863 | .9857 | .9846 | .9858 | .9825 | .9869 | .9869 | .9836 | .9665 | .9658 | .9484 | .9789 |
| 120 | .6297 | .6297 | .6280 | .6226 | .6184 | .6278 | .6301 | .6262 | .5962 | .6204 | .6256 | .6306 | .6211 | .6252 | .6169 | .6220 | .5931 | .5945 | .5883 | .6335 | |
| 121 | .7405 | .6646 | .7430 | .7269 | .7238 | .7372 | .7419 | .7405 | - | .7358 | .7360 | .7377 | .7286 | .7239 | - | .7321 | .7044 | .3791 | .7301 | .5343 | .7364 |
| 122 | .8497 | .8372 | .8483 | .8553 | .8566 | .8434 | .8480 | .8405 | .8275 | .8473 | .8472 | .8407 | .8468 | .8415 | .8439 | .8512 | .8476 | .8361 | .8289 | .8382 | .8394 |
| 123 | .8321 | .8351 | .8357 | .8362 | .8451 | .8348 | .8369 | .8325 | .8134 | .8396 | .8387 | .8263 | .8393 | .8356 | .8391 | .8321 | .8454 | .8307 | .7974 | .8278 | .8372 |
| 124 | .8367 | .8291 | .8420 | .8501 | .8500 | .8356 | .8483 | .8476 | .8155 | .8333 | .8369 | .8376 | .8313 | .8413 | .8405 | .8442 | .8483 | .8150 | .8083 | .8304 | .8483 |
| 126 | .5933 | .5957 | .5622 | .6140 | .6288 | .5867 | .6032 | .6097 | .5865 | .5575 | .5913 | .5820 | .5755 | .5842 | .6053 | .6047 | .6120 | .6042 | .5850 | .5665 | .6103 |
| 127 | .9696 | .9698 | .9706 | .9794 | .9778 | .9648 | .9709 | .9630 | .9587 | .9732 | .9746 | .9664 | .9657 | .9752 | .9780 | .9702 | .9631 | .9642 | .9642 | .9634 |
| 128 | .9806 | .9751 | .9751 | .9849 | .9799 | .9709 | .9773 | .9703 | .9669 | .9782 | .9760 | .9717 | .9679 | .9690 | .9763 | .9856 | .9775 | .9451 | .9748 | .9642 | .9669 |
| 129 | .9941 | .9921 | .9886 | .9878 | .9858 | .9870 | .9724 | .9893 | .9799 | .9878 | .9795 | .9862 | .9917 | .9925 | .9913 | .9909 | .9866 | .9941 | .9882 | .9925 | .9783 |
| 134 | .9167 | .7958 | .9038 | .8772 | .8896 | .7810 | .8080 | .7791 | .6077 | .9132 | .7776 | .7036 | .6591 | .6715 | .7299 | .8851 | .7657 | .3213 | .8332 | .5439 | .7352 |
| 136 | .5077 | .5373 | .5280 | .6721 | .6649 | .5144 | .6277 | .5765 | .5104 | .5292 | .5764 | .4991 | .5115 | .5045 | - | .6663 | .5831 | .5052 | .5160 | .4811 | .5143 |
| 137 | .3176 | .3594 | .3697 | .3591 | .3579 | .3330 | .3555 | .3636 | .3473 | .3185 | .3200 | .3200 | .3327 | .3073 | .3421 | .3797 | .3479 | .3500 | .3000 | .3170 | .3206 |
| 139 | .9074 | .9049 | .9085 | .9070 | .9082 | .9028 | .9095 | .9050 | .9035 | .9029 | .9092 | .9021 | .8998 | .9016 | .9039 | .9071 | .9041 | .8887 | .8945 | .8988 | .9060 |
| 142 | .5219 | .5302 | .5404 | .4812 | .5185 | .5333 | .5476 | .5435 | .5331 | .5336 | .5362 | .5513 | .5505 | .5450 | .5297 | .5641 | .4732 | .5309 | .5564 | .5358 | .5513 |
| 143 | .9333 | .9336 | .9405 | .9386 | .9388 | .9284 | .9376 | .9358 | .9164 | .9403 | .9400 | .9214 | .9373 | .9408 | .9590 | .9291 | .9239 | .9306 | | |
| 145 | .9946 | .9931 | .9884 | .9865 | .9923 | .9879 | .9821 | .9798 | .7557 | - | .9651 | .9882 | .9830 | .9811 | .9958 | .9927 | .9516 | .8963 | .8902 | .9370 | .9553 |
| 146 | .9527 | .9430 | .9565 | .9557 | .9575 | .9345 | .9634 | .9479 | .9332 | .9425 | .9513 | .9312 | .9263 | .9418 | .9476 | .9550 | .9552 | .8820 | .8970 | .8983 | .9422 |
| 147 | .5690 | .5487 | .5819 | .5575 | .5702 | .5503 | .5612 | .5729 | .5503 | .5627 | .5727 | .5614 | .5505 | .5555 | .5623 | .5837 | .5862 | .5214 | .5186 | .5252 | .5006 |
| 150 | .9659 | .9635 | .9647 | .9675 | .9689 | .9687 | .9665 | .9684 | .9709 | .9648 | .9688 | .9699 | .9680 | .9681 | - | .9713 | .9700 | .9679 | .9677 | .9863 | .9677 |
| 151 | .8517 | .8309 | .8800 | .7675 | .7637 | .7248 | .7664 | .7359 | .6981 | .7448 | .7208 | .7239 | .7090 | .7008 | - | .7422 | .7380 | .6812 | .7447 | .6869 | .7508 |
| 154 | .5718 | .5930 | .5661 | .6043 | .6016 | .5688 | .4750 | .5953 | .6106 | .5910 | .5681 | .5609 | .5901 | .5749 | .5706 | .5824 | .5853 | .5211 | .5390 | .5503 | .5473 |
| 157 | .7741 | .7390 | .7755 | .7817 | .7796 | .7580 | .7588 | .7556 | .6225 | .7723 | .7605 | .7554 | .7546 | .7567 | .7577 | .7782 | .7795 | .7098 | .6787 | .7319 | .7550 |

```
159  .9617  .9559 .9611 .9596 .9631 .9606 .9628 .9582 .9557 .9607   -   .9584 .9586 .9578 .9584 .9603 .9421 .8288 .9343 .9390 .9577
160  .6396  .6385 .6265 .6436 .6462 .6277 .6387 .6410 .6391 .6039 .6379 .6242 .6215 .6313   -   .6426 .6334 .6091 .6128 .6152 .6419
163  .8243  .8239 .8237 .8260 .8244 .8241 .8255 .8253 .8249 .8254 .8251 .8259 .8238 .8236   -   .8270 .8236 .8140 .8068 .8157 .8249
164  .9493  .9485 .9482 .9502 .9482 .9449 .9471 .9467 .9444 .9469 .9460 .9469 .9460 .9448 .9455 .9485 .9489 .9431 .9467 .9408 .9461
166  .9868  .9837 .9885 .9864 .9866 .9848 .9876 .9839 .9841 .9856 .9858 .9849 .9830 .9854 .9906 .9841 .9841 .9841 .9832 .9841
167  .9565  .9557 .9459 .9677 .9652 .9448 .9350 .9502 .9079   -     -   .9439 .9603 .9531   -   .9669 .9488 .9624 .8934 .9344 .9130
168  .3915  .3989 .3926 .3859 .3844 .3851 .3986 .3692 .4027 .3779 .3945 .3590 .3726 .3920 .3912 .3949 .3935 .3926 .4058 .3972 .3880
169  .9327  .9285 .9274 .9210 .9218 .9246 .9253 .9249 .9107 .9242 .9241 .9224 .9224 .9249 .9229 .9114 .9148 .8999 .9233
170  .9911  .9847 .9819 .9449 .9568 .9569 .9535 .9564 .8195 .9542 .9367 .9456 .9522 .9500 .9465 .9460 .8841 .6248 .9222 .6516 .9208
171  .9612  .9125 .9652 .9246 .9169 .8550 .8751 .8555 .8318 .8680 .8744 .8518 .8465 .8325   -   .9238 .8577 .7527 .8433 .7804 .8554
173  .9828  .9696 .9776 .9474 .9469 .9494 .9713 .9519 .7543 .9719 .9683 .9537 .9395 .9586 .9505 .9533 .9346 .7634 .9031 .7879 .9504
174  .9941  .9250 .9815 .7220 .7188 .6122 .7159 .6287 .6082 .7372 .5958 .5930 .5981 .6201   -   .7246 .6449 .4977 .5782 .5392 .6015
175  .9088  .6140 .6595 .6325 .6152 .5708 .5945 .5741 .5611 .6330   -   .5745 .5760 .5706 .5965 .6393 .6057 .5558 .5802 .5582 .5745
177  .5995  .5821 .5891 .6123 .6111 .5763 .5751 .5712 .4357 .5648 .5524 .5710 .5822 .5635   -   .6011 .6166 .4845 .5989 .5071 .5718
180  .9949  .9954 .9948 .9942 .9949 .9941 .9920 .9940 .9880 .9957 .9924 .9935 .9941 .9915   -   .9947 .9940 .9896 .9946 .7259 .9927
182  .9211  .9380 .9295 .9235 .9381 .9012 .8606 .9147 .9131 .8937   -   .8980 .9049 .9154   -   .9361 .9038 .8646 .7853 .8919 .8924
183  .9129  .8913 .9172 .9132 .9105 .8938 .9288 .8992 .9184 .9078 .9294 .9279 .9221 .9041 .8447 .9068 .9169
185  .7213  .7170 .7217 .7233 .7230 .7181 .7217 .7087 .7205 .7240 .7245 .7194 .7175 .7171 .7216 .7221 .7198 .7125 .7145 .7188 .7174
186  .7811  .7122 .7004 .5778 .5166 .6549 .5081 .5953 .4919 .5344   -   .6126 .6532 .5874 .4925 .6027 .5498 .7868 .5556 .5191 .5339
187
193  .9800    -   .9785 .9789 .9794 .9783 .9789 .9786 .9754 .9794 .9793 .9780 .9790 .9789 .9793 .9781 .9803 .9798 .9799 .9780 .9790
194  1.000  .9766 1.000 1.000 1.000 .9574 1.000 .9780 .8449 1.000 1.000 .9565 .9694 .9760 .9841 1.000 1.000 1.000 .8810 .9381 .9460
195  .7590  .7574 .7406 .7611 .7622 .7082 .7570 .7490 .7225   -   .7096 .7163 .7135   -   .7626 .7441 .6208 .6882 .6338 .7453
196  .9292  .9110 .9309 .9298 .9308 .8012 .7973 .8002 .7976 .8045 .7980 .7984 .7974 .7982   -   .9309 .9311 .7954 .9310 .7973 .7975
197  .5284  .5244 .5102 .5362 .5358 .5084 .5348 .5157 .4508   -   .4737 .4693 .4859 .5014 .5270 .4729 .4169 .4065 .4351 .3504
198  .8090  .8059 .8165 .8144 .8084 .8125 .8141 .8109 .8131 .8128 .8078 .8130 .8136 .8114 .8089 .8111 .8129 .8167 .7984 .8110
199  .9950  .9248 .9898 .8610 .8647 .9527 .9761 .9674 .8551 .9907 .9734 .9211 .8722 .8560 .9607 .8659 .8262 .6801 .8814 .7515 .9517
200  .8551  .8689 .8474 .8627 .8624 .8581 .8534 .8553 .8502 .8526 .8471 .8540 .8575 .8648 .8616 .8744 .8698 .8507 .8602 .8575 .8494
201  .8796  .8669 .8760 .8836 .8774 .8496 .8487 .8507 .8328 .8787 .8544 .8489 .8475 .8464 .8509 .8870 .8827 .8581 .8372 .8420 .8474
203  .9388  .7733 .9156 .6044 .7900 .8606 .9138 .8889 .5547 .8866 .9089 .8570 .7952 .7808 .8821 .8959 .7083 .2794 .7277 .3841 .8446
204  .9371  .7835 .9038 .8579 .8134 .8727 .9106 .9005 .5530   -   .8706 .8567 .7966 .7695 .8762 .8770 .7101 .2800 .7153 .3822 .8453
206  .7318  .7348 .7360 .7292 .7349 .7343 .7331 .7259 .7349   -   .7384 .7267 .7372 .7376 .7350 .7343 .7326 .7328 .7203 .7279 .7159
207  .7388  .7361 .7370 .7435 .7431 .6068 .7361 .7348 .7393   -   .7344 .7342 .7380 .7312 .7374 .7355 .7450 .7203 .7234 .7332 .7298
208  .9854  .9766 .9849 .9643 .9709 .9745 .9759 .9722 .8842 .9739 .9706 .9652 .9724 .9597   -   .9650 .9102 .6940 .9508 .8760 .9740
209  .8332  .6892 .7240 .8632 .8643 .6152 .6691 .5917 .7645 .6248 .5856 .5818 .6226   -   .8638 .7412 .6035 .6783 .6066 .5729
210  .9861  .9847 .9872 .9867 .9883 .9867 .9871 .9865 .9821 .9871 .9864 .9878 .9861 .9850 .9872 .9876 .9852 .9819 .9866 .9707 .9865
211  .8289  .7993 .8325 .8213 .8128 .7241 .7123 .6677 .7392 .8125 .7360 .7291 .7097 .7015 .7149 .8332 .8204 .6256 .8276 .6158 .6939
213  .9722  .9713 .9693 .9660 .9680 .9615 .9642 .9577 .9383 .9637 .9642 .9647 .9660   -   .9688 .9655 .9775 .9625 .9620 .9593
214  .8648  .8408 .8653 .8572 .8718 .8445 .8523 .8467 .8330 .8523 .8555 .8407 .8483 .8438 .8578 .8623 .8602 .8350 .8650 .8340 .8457
215  .9653  .9690 .9557 .9400 .9595 .9563 .9555 .9535 .9377 .9575 .9495 .9542 .9610 .9503 .9468 .9502 .9472 .9425 .9600 .9498 .9512
216  .7553  .7535 .7633 .7488 .7428 .7643 .7668 .7592 .6900 .7615 .7445 .7762 .7580 .7662 .7650 .7533 .7490 .7492 .7600 .6218 .7577
217  .9650  .9623 .9667 .9610 .9520 .9623 .9490 .9600 .9455   -   .9588 .9642 .9592   -   .9625 .9580 .9500 .9600 .9572 .9505
218  .7902  .8055 .8727 .7807 .7855 .8477 .8333 .8298 .6793 .8218 .7977 .8347 .8397 .8385 .8180 .7945 .7603 .8238 .8050 .8257 .8332
219  .9809  .9917 .9988 .9630 .9907 .9923 .9802 .9948 .9880 .9812 .9793 .9914 .9929 .9861 .9796 .9861 1.000 .9772 .9864
220  .6308  .6162 .6328 .6203 .6189 .6300 .6276 .6290 .6116 .6277 .6259 .6220 .6190 .6284   -   .6192 .6112 .5680 .5863 .5964 .6263
221  .9843  .9600 .9893 .9805 .9803 .9832 .9817 .9826 .9613 .9833 .9814 .9807 .9786 .9767   -   .9795 .9701 .8369 .9668 .9034 .9798
222  .9428  .9124 .9426 .9459 .9466 .9194 .9276 .9250 .9273 .9386 .9279 .9148 .9252 .9325 .9458 .9456 .7785 .9112 .8502 .9185
224  .9690  .9704 .9699 .9695 .9695 .9697 .9691 .9681 .9619   -   .9625 .9694 .9697 .9699   -   .9694 .9638 .9583 .9656 .9635 .9682
225  .9989  .9838 .9969 1.000 1.000 .9809 .9965 .9845 .9979 .9998 .9773 .9730 .9893 1.000 1.000 .9964 .8788 .9504 .9822
226  .7475  .7472 .7322 .7514 .7508 .7338 .7500 .7481 .7309 .7312 .7407 .7323 .7335 .7328   -   .7464 .7407 .7237 .7249 .7295 .7459
227  .8408  .7850 .7898 .7715 .8533 .8240 .8115 .8060 .7035 .8125 .7360 .7717 .8038 .7850 .7804 .7640 .7804 .6531 .7781 .7379 .7167
228  .7054  .7215 .7294 .5958 .6510 .6981 .6921 .6938 .2073 .7077 .6335 .6310 .6606 .6927 .5781 .6269 .4640 .6250 .3869 .6348
229  .8752  .8456 .8583 .8033 .8896 .8717 .8473 .8531 .7575 .8740 .8140 .8123 .8388 .8550 .8260 .8125 .7952 .8125 .8312 .7979 .7917
230  .9053  .8898 .8964 .8995 .9063 .8939 .9019 .9032 .8875 .8894 .8998 .8958 .8958 .8916 .9016 .9012 .9038 .8824 .8739 .8905 .8989
231  .9913  .9867 .9824 .9824 .9825 .9779 .9571 .9785 .9698 .9805 .9868 .9811   -   .9831 .9801 .9582 .9795 .9810 .9807
232  .9466  .9437 .9442 .9479 .9440 .9374 .9483 .9431 .9369 .9337 .9452 .9303 .9474 .9466 .9408 .9481 .9500 .9408 .9507 .9439 .9446
234  .9668  .9640 .9677 .9745 .9667 .9669 .9671 .9663 .9295 .9685 .9632 .9694 .9679 .9646 .9707 .9712 .9703 .9565 .9616 .8715 .9661
235  .9330  .9291 .9336 .9381 .9351 .9273 .9309 .9183 .9324 .9339 .9291 .9381 .9375 .9357 .9324 .9369 .9258 .9318
236  .8937  .8816 .8846 .8856 .8880 .8820 .8978 .8867 .8978 .8824 .8969 .8907 .8920 .8999 .8912 .8773 .8941 .8871 .8914 .8835 .8909
237  .8943  .9048 .8989 .9151 .9016 .8989 .8959 .8975 .8779 .8986 .8945 .9005 .8989 .9021 .8957 .9110 .8913 .9027 .8904 .8854 .8995
238  .9940  .9938 .9943 .9414 .9913 .9950 .9935 .9942 .9807 .9929 .9937 .9943 .9509 .9879 .9281 .9909 .9854 .9927
239  .8860  .8771 .8892 .8766 .8797 .8684 .8809 .8710 .8538 .8698 .8654 .8610 .8583 .8480 .8744 .8807 .8852 .7294 .8686 .7987 .8632
240  .9885  .9865 .9887 .9824 .9860 .9854 .9914 .9882 .9808 .9904 .9853 .9853 .9845 .9881 .9878 .9818 .9732 .8815 .9534 .9569 .9841
243  .7345  .8043 .7452 .7556 .7782 .7291 .7547 .7585 .7492 .7205 .7532 .7305 .7333 .7436 .7401 .7587 .7628 .6973 .7006 .7217 .7473
246  .8635  .8294 .8632 .8645 .8613 .8932 .8537 .8720 .6768 .8537 .7953 .8926 .8970 .8831 .8645 .8616 .8515 .8711 .8673 .8657 .8720
249  .9268  .9283 .9216 .9189 .9223 .9179 .9204 .9192 .9191 .9175 .9192 .9163 .9203 .9198 .9179 .9234 .9218 .9213 .9239 .9178 .7902
250  .9802    -   .9795 .9715 .9711 .9741 .9765 .9696 .9688 .9779 .9774 .9726 .9735 .9625 .9741 .9367 .7705 .7757 .9221 .9746
251  .8375  .7899 .8824 .7922 .7858 .6600 .7196 .6488 .6475 .7416 .7168 .6581 .6596 .6601 .6840 .7983 .7490 .5718 .6640 .6227 .6429
253  .9104  .9051 .9131 .9142 .9135 .8993 .8959 .9045 .8612 .9040 .8947 .8979 .8923 .8868 .9125 .9162 .9086 .8367 .9075 .8654 .8953
254  .9277  .9042 .9190 .9149 .9202 .9199 .9137 .9188 .7957 .8986 .9059 .9185 .9130 .9110 .9177 .9196 .8297 .8896 .8420 .9052
255  .9304  .9338 .9262 .9313 .9293 .9277 .9280 .9250 .9340 .9219 .9282 .9215 .9251 .9333 .9224 .9295 .9340 .9295 .9246 .9280 .9297
256  .9220  .9457 .9168 .9131 .9204 .9168 .9080 .9250 .8729 .8978 .9080 .9166 .9285 .9289   -   .9166 .9154 .8621 .9028 .8966 .9066
260  .9996  .9992 .9994 .9998 .9997 .9991 .9996 .9989 .9997 .9995 .9988 .9992   -   .8539 .9997 .9268 .9991 .9750 .9990
263  .9320  .9404 .9388 .9523 .9467 .9386 .9393 .9383 .9399 .9304 .9346 .9374 .9397 .9373 .9407 .9461 .9330 .9191 .9175 .9301 .9389
264  .9437  .9329 .8310 .9558 .9529 .8735 .9493 .9456 .7730   -     -   .8410 .8787 .9078 .9475 .9591 .9428 .7972 .8480 .8504 .9354
265  .8377  .8313 .8407 .8340 .8303 .8330 .8223 .7017 .8163 .8053 .8187 .8353 .8181 .8250 .8217 .8350
266  .7250  .7200 .7297 .7467 .7430 .7250 .7123 .7123 .7227 .7187 .7303 .7273 .7307 .7087 .7130 .7433 .7450 .7350 .7250 .7230 .7113
268  .7512  .7482 .7558 .7794 .7911 .7536 .7632 .6769 .5748 .7606 .7539 .7326 .7465 .7477 .7741 .7805 .7654 .7253 .7532 .4896 .7424
273  .9588  .9476 .9638 .9454 .9524 .9309 .9471 .9357 .9154 .9463 .9330 .9287 .9337 .9410 .9424 .9316 .9161 .8820 .9112 .9315
274  .9674  .9695 .9678 .9702 .9718 .9685 .9674 .9662 .9678 .9662 .9705 .9686 .9689 .9672   -   .9712 .9717 .9677 .9714 .9522 .9670
275  .7969  .7976 .7973 .7962 .8057 .7925 .8005 .8002 .7988 .7931 .8005 .7913 .7914 .7906 .7917 .8031 .7968 .7944 .7793 .7840 .7977
276  .9974  .9996 .9981 .9784 .9900 .9964 .9962 .9967 .8820 .9962 .9927 .9953 .9985 .9968 .9967 .9081 .9675 .9982 .9855 .9800 .9932
277  .9921  .9815 .9896 .9952 .9949 .9867 .9894 .9840 .9756 .9910 .9915 .9842 .9840 .9876 .9883 .9953 .9964 .9382 .9500 .9745 .9811
278  .9873  .9783 .9861 .9915 .9952 .9842 .9858 .9807 .9353 .9871 .9858 .9814 .9830 .9814 .9847 .9974 .9929 .9571 .9563 .9727 .9759
279  .6692  .6889 .6814 .6921 .6924 .6956 .6932 .6924 .6350 .6870 .6938 .6952 .6904 .6921 .6852 .6858 .6929 .6971 .6607 .6836 .6927
282  .5243  .4937 .5102 .5103 .5170 .5026 .5153 .5070 .4751 .4803 .5033 .4988 .4934 .5023 .5151 .5147 .5090 .4357 .4777 .4660 .5139
283  .9778  .9799 .9762 .9763 .9774 .9773 .9782 .9676 .9791 .9782 .9740 .9783 .9748 .9796 .9715 .9757 .9745 .9791 .9797 .9728 .9762
284  .8024  .7796 .8576 .7604 .7529 .8235 .7824 .8439 .7820 .8133 .7784 .8043 .7945 .7757 .8024 .7812 .7278 .7596 .7294 .7165 .7976
285  .7718  .7384   -   .6974 .7088 .6985 .7184 .7005 .6393   -   .6870 .6910 .6869   -   .6960 .6271 .5742 .6741 .6341 .6992
287  .6774  .6449 .6713 .6519 .6617 .6434 .6537 .6457   -   .6500 .6519 .6412 .6169 .6143   -   .6519 .6159 .2455 .6310 .3154 .6433
288  .9877  .9349 .9773 .9958 .9965 .9341 .9797 .9282 .9096 .9889 .9839 .9274 .9189 .9253 .9606 .9982 .9930 .6868 .9249 .8440 .9205
289  .8990  .8931 .9050 .9025 .8998 .9078 .9013 .9056 .9006 .9072 .8978 .9072 .9049 .8998 .9088 .9041 .9001 .8950 .9056 .9032 .9040
```

| ID | | | | | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 290 | .8624 | .8659 | .8604 | .8541 | .8596 | .8557 | .8611 | .8588 | .8625 | .8624 | .8609 | .8602 | .8613 | .8609 | .8498 | .8495 | .8510 | .8628 | .8510 | .8637 | .8119 |
| 291 | .3283 | .3435 | .3329 | .3392 | .3391 | .3351 | .3330 | .3355 | .3346 | .3311 | .3326 | .3349 | .3390 | .3384 | .3374 | .3481 | .3372 | .3367 | .3560 | .3377 | .3338 |
| 292 | .8702 | .8522 | .8626 | .8369 | .8438 | .8587 | .8642 | .8561 | .8569 | .8551 | .8504 | .8569 | .8565 | .8564 | .8499 | .8173 | .8291 | .8620 | .8410 | .8590 | .8587 |
| 294 | .8753 | .8745 | .8972 | .8878 | .9043 | .8954 | .9014 | .9087 | .8962 | - | - | .8859 | .8728 | .8836 | .8819 | .8898 | .8802 | .8450 | .8708 | .8578 | .8145 |
| 296 | .7448 | .7435 | .7393 | .7500 | .7534 | .7268 | .7234 | .7211 | .7517 | .7298 | .7220 | .7190 | .7280 | .7198 | .7260 | .7339 | .7386 | .7045 | .7710 | .6881 | .7212 |
| 298 | .6363 | .6119 | .6127 | .6362 | .6073 | .5867 | .5819 | .5617 | .5823 | .6065 | .5515 | .5017 | .5685 | .5746 | .5550 | .6360 | .6329 | .5481 | .6562 | .4879 | .5648 |
| 299 | .6327 | .5924 | .6265 | .6242 | .6315 | .5779 | .5517 | .5586 | .5430 | .5942 | .5327 | .5498 | .5752 | .5439 | .5816 | .6295 | .6235 | .5306 | .6337 | .5235 | .5531 |
| 300 | .6034 | .6305 | .5984 | .5984 | .6083 | .6031 | .5982 | .5928 | .5957 | .5987 | .5672 | .5912 | .5915 | .6002 | .5138 | .5825 | .6247 | .5859 | .5960 | .5437 | .5908 |

Table 11: Detailed results (average RMSE) for all regression datasets over all methods. Due to space constraints, dataset IDs from Ye et al. (2024a) are used in place of dataset names."XGB", "RF", "LG", "RN", "ND", "ST", "TG", "DAN", "FTT", "DCN", "TT", "PT", "EF", and "TR" denote "XGBoost", "Random Forest", "LightGBM", "ResNet","NODE", "SwitchTab", "TANGOS", "DANets", "FT-T", "DCNv2", "PTaRL","Excelformer" and "TabR", respectively.

| ID | M-NCA | L-NCA | TR | XGB | CGB | MLP | FTT | DCN | ND | MP | EF | DAN | TG | ST | LG | RF | SVM | KNN | SW | PT |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 $(\times 10^3)$ | .6383 | .4786 | .6786 | .4711 | .4368 | .7012 | .6750 | .5312 | .5102 | .6531 | .6543 | 1.661 | .6043 | .6453 | .6143 | .4377 | .7804 | .4864 | .8309 | .6057 |
| 2 $(\times 10)$ | .1002 | .1007 | .1011 | .1009 | .1001 | .1008 | .1005 | .1003 | .1007 | .1002 | .1004 | .4977 | .1020 | - | .1002 | .1020 | .2380 | .1154 | .4332 | .1011 |
| 3 $(\times 10^{-5})$ | 0.000 | 437.8 | 20.88 | 219.7 | 33.34 | 7578. | 1685. | 4921. | Inf | 123.5 | 6205. | Inf | Inf | 5003. | 1531. | Inf | Inf | 0.000 | Inf | 9680. |
| 4 $(\times 10^{-5})$ | 0.000 | 251.3 | 16.08 | 111.0 | 3.716 | 7296. | 4760. | 5418. | Inf | 907.7 | 6965. | Inf | Inf | 990.8 | 17.89 | 8267. | Inf | 0.000 | Inf | Inf |
| 6 $(\times 10)$ | .2106 | .2123 | .2101 | .2167 | .2179 | .2130 | .2108 | .2097 | .2225 | .2134 | .2182 | .3206 | .2136 | .2125 | .2173 | .2155 | .2390 | .2284 | .3245 | .2102 |
| 7 $(\times 10^{-3})$ | .1576 | .1560 | .1522 | .1527 | .1465 | .1563 | .1554 | .2719 | .1519 | .1540 | .1528 | .3841 | .1564 | - | .1503 | .1559 | .1696 | .2047 | .3832 | .1554 |
| 9 $(\times 10^3)$ | .7892 | .7739 | .7909 | .7490 | .8173 | .7743 | .7743 | .7869 | .7601 | .7708 | .7708 | 2.326 | .7859 | .7713 | .7473 | .7509 | .7617 | .7465 | 1.826 | .7617 |
| 13 $(\times 10^2)$ | .1143 | .1152 | .1142 | .1149 | .1143 | .1173 | .1153 | .1171 | .1150 | .1161 | .1147 | .1555 | .1175 | .1159 | .1143 | .1141 | .1188 | .1158 | .1564 | .1166 |
| 14 $(\times 10^3)$ | .4549 | .4593 | .4594 | .4560 | .4557 | .4631 | .4547 | .4583 | .4556 | .4592 | .4562 | .7285 | .4626 | - | .4555 | .4617 | .4855 | .4620 | .7180 | .4572 |
| 15 $(\times 10)$ | .4574 | .4622 | .4599 | .4583 | .4673 | .4603 | .4685 | .4619 | .4640 | .4602 | 1.099 | .4709 | .4699 | - | .4569 | .4581 | .6324 | .4904 | 1.026 | .4650 |
| 16 $(\times 10)$ | .2892 | .3072 | .2933 | .2922 | .2902 | .3036 | .2930 | .3015 | .2909 | .2961 | .2925 | .6399 | .3044 | - | .2953 | .2988 | .4536 | .3112 | .6467 | .3006 |
| 20 | .6006 | .6404 | .5728 | .8889 | .7607 | .7639 | .7017 | .7239 | 1.161 | .7646 | .8134 | 2.978 | .7849 | .7574 | .9117 | 1.053 | 1.518 | 1.082 | 2.982 | .7976 |
| 21 | .5555 | .6505 | .4548 | .6484 | .5369 | .6013 | .5398 | .5542 | .8019 | .5597 | .5597 | 2.496 | .5998 | .5565 | .6473 | .7982 | 1.038 | .7855 | 2.407 | .6255 |
| 22 $(\times 10^{-1})$ | .4420 | .1828 | .3426 | .6955 | .3372 | .2649 | .1741 | 1.289 | .3576 | .3547 | .2355 | 7.373 | .6013 | .3766 | .5384 | .5066 | 3.245 | 1.300 | 7.519 | .7058 |
| 24 $(\times 10^3)$ | .4618 | .4782 | .5150 | .4526 | .4873 | .5117 | .5602 | .5152 | .5497 | .4886 | .5139 | 1.776 | .5061 | .5379 | .4650 | .4667 | 1.530 | .6228 | 1.791 | .5148 |
| 25 $(\times 10^4)$ | .1245 | .1317 | .1458 | .1199 | .1275 | .1408 | .1382 | .1330 | .1401 | .1363 | .1308 | .8299 | .1316 | .1355 | .1256 | .1228 | .1552 | .1347 | .1748 | .1326 |
| 26 $(\times 10^3)$ | .3655 | .3678 | .3664 | .3654 | .3646 | .3659 | .3655 | .3650 | .3642 | .3655 | .3658 | .4362 | .3655 | .3658 | .3646 | .3692 | .3842 | .3671 | .4395 | .3638 |
| 30 | .2011 | .2293 | .2427 | .2749 | .2638 | .2344 | .2768 | .2360 | .3040 | .2252 | .2697 | .4986 | .2148 | .2568 | .2769 | .3052 | .3448 | .3293 | .5014 | .2272 |
| 31 | .2342 | .2259 | .2262 | .2691 | .2592 | .2228 | .2625 | .2171 | .3034 | .2243 | .2603 | 2.297 | .2100 | .2587 | .2711 | .2856 | .3459 | .3148 | .5101 | .2129 |
| 32 | .2445 | .2198 | .2368 | .2738 | .2674 | .2433 | .2839 | .2312 | .2315 | .2444 | .2237 | | .2892 | .2795 | .2945 | .3564 | .3324 | | .5614 | .2381 |
| 33 | .2449 | .2285 | .2433 | .2829 | .2721 | .2701 | .2815 | .2519 | .3080 | .2921 | .2836 | .5876 | .2465 | .2883 | .2868 | .3046 | .3633 | .3389 | .5202 | .2514 |
| 34 | .2483 | .2460 | .2796 | .3098 | .2954 | .2609 | .2958 | .2645 | .3271 | .2792 | .2913 | .8969 | .2355 | .2940 | .3105 | .3255 | .3602 | .3371 | .5456 | .2456 |
| 35 $(\times 10)$ | .1487 | .1507 | .1549 | .1483 | .1486 | .1504 | .1520 | .1511 | .1497 | .1508 | .1512 | .1594 | .1498 | .1502 | .1483 | .1488 | .1532 | .1491 | .1567 | .1496 |
| 38 $(\times 10^5)$ | .4922 | .5366 | .6506 | .6012 | .6000 | .6548 | .6548 | .6771 | .6549 | .6539 | .6540 | .6108 | .6557 | .6554 | .5992 | .6011 | .6327 | .6117 | .6288 | .6133 |
| 41 $(\times 10^3)$ | .5250 | .5269 | .5116 | .5328 | .5249 | 1.228 | .5328 | 24.39 | .5882 | .5517 | .5320 | 4.049 | .9922 | - | .5322 | .5547 | 1.507 | .7572 | 3.918 | .6726 |
| 54 $(\times 10^2)$ | .1735 | .2273 | .1677 | .1782 | .1588 | .2499 | .1794 | .4945 | .1877 | .1779 | .1803 | .3016 | - | .1777 | .1700 | .1690 | .2841 | .2399 | .2958 | .4594 |
| 57 $(\times 10^2)$ | .7280 | .7885 | .7423 | .7410 | .7363 | .7891 | .7866 | .7890 | .7855 | .7453 | .7866 | .9422 | .7893 | - | .7405 | .7633 | .8490 | .8072 | .9358 | .7869 |
| 62 | .4688 | .4866 | .4839 | .4320 | .4320 | .4844 | .4843 | .6600 | .4840 | .4760 | .4851 | .4889 | .4846 | .4854 | .4247 | .4287 | .4878 | .4687 | .4898 | .4845 |
| 64 | .2510 | .2598 | .3180 | .2628 | .2609 | .3121 | .2933 | .3307 | .3014 | .3015 | .3096 | .6651 | .2889 | .2943 | .2758 | .2720 | .3584 | .2514 | .5154 | .2924 |
| 66 $(\times 10^{-1})$ | .2629 | .4044 | .1991 | .4063 | .3483 | .2844 | .2959 | .2605 | .5025 | .2647 | .3396 | .2896 | .2816 | .4200 | .4478 | .5894 | .8366 | 3.200 | .3248 | |
| 69 | .4293 | .4292 | .4331 | .4286 | .4288 | .4592 | .4397 | .4827 | .4287 | .4669 | .4343 | .4480 | .4326 | .4483 | .4286 | .4301 | .4299 | .4303 | .4337 | .4286 |
| 70 | .2011 | .2007 | .2070 | .2218 | .2135 | .2369 | .2137 | .2327 | .2055 | .2209 | .2134 | .7146 | .2191 | .2187 | .2198 | .2293 | .2796 | .2345 | .5008 | .2030 |
| 71 | .3148 | .3134 | .3237 | .3075 | .3065 | .3280 | .3394 | .3294 | .3120 | .3320 | .3320 | .3186 | .3280 | .3215 | .3222 | .3060 | .3035 | .3203 | .3088 | .3299 |
| 73 $(\times 10)$ | 1.063 | 1.324 | 1.143 | 1.053 | 1.024 | 1.325 | 1.321 | 1.321 | 1.323 | 1.067 | 1.324 | 1.330 | 1.325 | - | .7682 | .7664 | 1.326 | 1.321 | 1.331 | 1.322 |
| 74 | .7224 | - | .9062 | .7055 | .7070 | .9059 | .7158 | .8344 | .7371 | .8890 | .7193 | 1.140 | .8121 | .8318 | .7078 | .7168 | .8272 | .8311 | 1.036 | .7902 |
| 76 $(\times 10^2)$ | .7057 | .7131 | .7350 | .7180 | .7264 | .7806 | .7633 | .7226 | .7830 | .7835 | .7290 | 1.864 | .7870 | .7159 | .7601 | 1.561 | 1.116 | 1.818 | | .7136 |
| 77 $(\times 10^3)$ | .4135 | .4026 | .4606 | .4448 | .4506 | .4789 | .4609 | .4652 | .5112 | .4887 | .4169 | .8800 | .4641 | .4705 | .4528 | .4587 | .6939 | .4968 | .8839 | .4777 |
| 78 $(\times 10^5)$ | .1563 | .1855 | .2318 | .1821 | .1870 | .2211 | .1922 | .2151 | .2601 | .2309 | .2030 | 3.268 | .2304 | - | .1803 | .1722 | .3027 | .8028 | 2.086 | .2251 |
| 79 $(\times 10^4)$ | .6636 | .6932 | .7048 | .6917 | .7173 | .8200 | .7724 | .7958 | .7568 | .7601 | .7965 | 12.36 | .8391 | - | .7961 | .7294 | 1.384 | 3.229 | 12.11 | .8499 |
| 80 $(\times 10^5)$ | .2513 | .4257 | .2175 | .1793 | .1878 | .2345 | .1909 | .3327 | .3006 | .1890 | .2061 | .2687 | - | .1944 | .1804 | .3290 | .4312 | 1.050 | .2356 | |
| 81 $(\times 10^4)$ | .3750 | .6091 | .3894 | .4923 | .4498 | .5998 | .4370 | .7411 | .5431 | .4218 | .4144 | 10.12 | .5898 | - | .4926 | .6289 | 1.465 | 1.176 | 7.202 | .5841 |
| 82 $(\times 10^{-1})$ | .9290 | .8924 | 1.010 | .9142 | 1.221 | 1.233 | .8820 | 1.332 | 2.946 | .8418 | .8758 | 12.42 | 1.300 | 1.259 | .7684 | .9063 | 3.296 | 2.819 | 4.715 | 1.133 |
| 86 $(\times 10^5)$ | .8508 | .9156 | .8724 | .8740 | .8034 | .9007 | .8806 | .8820 | .9576 | .8993 | .8537 | 3.174 | .9223 | .9105 | .8812 | .9800 | 2.050 | 1.014 | 3.184 | .9527 |
| 87 $(\times 10^3)$ | .7965 | .8100 | .8886 | .7033 | .7749 | 1.112 | .8092 | .8198 | 1.372 | .9353 | - | 2.959 | 1.107 | .7871 | .7970 | .7753 | 2.312 | 1.113 | 2.958 | .7946 |
| 89 $(\times 10^2)$ | .4025 | .4025 | .4031 | .4069 | .4037 | .4066 | .4028 | .4052 | .4028 | .4034 | .4036 | .6632 | .4064 | - | .4052 | .4097 | .4488 | .4141 | .7061 | .4053 |
| 90 $(\times 10^2)$ | .1642 | .1651 | .1563 | .1559 | .1582 | .1552 | .1530 | .1587 | .1531 | .1530 | .1541 | .1963 | .1549 | .1532 | .1564 | .1560 | .1750 | .1608 | .1981 | .1553 |
| 92 $(\times 10^2)$ | .1552 | .1561 | .1825 | .1557 | .1548 | .1724 | .1782 | .1546 | .1606 | .1705 | .1739 | .1948 | .1545 | .1763 | .1556 | .1561 | .1644 | .1586 | .1856 | .1649 |
| 93 $(\times 10)$ | .5423 | .4981 | .1515 | .1599 | .8308 | .8344 | .8324 | .8333 | .8336 | .8346 | .8333 | .8360 | .1502 | .1710 | .7772 | .4835 | .5872 | .3057 | | |
| 94 | .4283 | .4273 | .4439 | .4294 | .4425 | .4691 | .4388 | .4516 | .4465 | .4629 | .4418 | .6890 | .4806 | .4886 | .4334 | .4302 | .4264 | .4271 | .5081 | .4285 |
| 96 $(\times 10)$ | .3083 | .3320 | .3047 | .3661 | .3475 | .3675 | .3649 | .3549 | .4230 | .3593 | .3577 | .6248 | .3530 | - | .3744 | .3911 | .5237 | .3690 | .6064 | .3474 |
| 105 | .4245 | .4277 | .5212 | .4262 | .4266 | .5028 | .4364 | .4544 | .5548 | .4287 | .5446 | .4500 | .4261 | .4292 | .4551 | .4428 | .4979 | .4298 | | |
| 106 $(\times 10^2)$ | .2903 | .2707 | .2875 | .2863 | .2864 | .2874 | .2868 | .2884 | .2873 | .2867 | .2861 | .2923 | .2869 | .2857 | .2848 | .2936 | .2859 | .2854 | .2889 | .2862 |
| 108 $(\times 10)$ | .4414 | .4475 | .4468 | .4047 | .4023 | .4454 | .4492 | .4613 | .4470 | .4446 | .4467 | .4402 | .4494 | .4507 | .4022 | .4087 | .4174 | .4112 | .4280 | .4124 |
| 109 $(\times 10)$ | .1923 | .2250 | .1929 | .1535 | .1849 | .1996 | .1488 | .2136 | .1648 | .1761 | .2081 | .4692 | .1859 | .1449 | .1270 | .1572 | .3060 | .4675 | .2009 | |
| 110 $(\times 10)$ | 1.038 | 1.034 | 1.032 | .9958 | .9911 | 1.044 | 1.065 | 1.030 | 1.269 | 1.030 | 1.061 | 3.713 | 1.088 | - | 1.022 | 1.057 | 1.838 | 1.078 | 3.361 | 1.079 |
| 113 $(\times 10)$ | .5990 | .8794 | - | 1.047 | .9139 | .9958 | 1.265 | 1.929 | .8625 | .6186 | - | 12.59 | .9678 | - | .9796 | 2.319 | 3.120 | .9004 | 12.02 | 1.044 |
| 114 $(\times 10^{-1})$ | 1.036 | 1.076 | 1.024 | 1.010 | 1.044 | 1.011 | 1.023 | 1.041 | 1.011 | 1.032 | 1.014 | 1.023 | 1.015 | 1.014 | 1.023 | 1.059 | .9969 | 1.025 | .9904 | |
| 116 | .4668 | .4737 | .5101 | .4589 | .4705 | .4595 | .4563 | .4681 | .4758 | .4777 | .4642 | .9514 | .4546 | .4747 | .4793 | .4565 | .5822 | .4519 | .5254 | .4413 |
| 118 | .6260 | .5809 | .6745 | .6192 | .6081 | .6693 | .6514 | .6754 | .6512 | .6669 | .6450 | .8191 | .6644 | .6637 | .6193 | .6183 | .6900 | .6532 | .8179 | .6598 |
| 119 | .6110 | .6137 | .6632 | .6235 | .6132 | .6925 | .7003 | .7280 | .6924 | .6853 | .7047 | .9056 | .6934 | .7077 | .6316 | .6466 | .7605 | .6277 | .9103 | .6933 |
| 125 $(\times 10)$ | .1613 | .6764 | .1260 | .1775 | .1390 | .1536 | .1262 | .1334 | .1302 | .1884 | .1576 | 1.386 | .1564 | .1431 | .1784 | .2015 | .4865 | .2720 | .7085 | .1565 |
| 130 $(\times 10^{-1})$ | .9463 | .9082 | .8969 | .8078 | .7636 | .9538 | .9164 | .9425 | 1.014 | .9466 | .9092 | 5.365 | 1.072 | .9971 | .8524 | .8655 | 4.616 | 1.097 | 5.204 | .9583 |
| 131 $(\times 10^3)$ | .3510 | .3949 | .3969 | .3332 | .3336 | .3705 | .3907 | .3986 | .3448 | .4057 | .3812 | .4632 | .3739 | .3915 | .3393 | .3454 | .4156 | .3557 | .4598 | .3611 |
| 132 $(\times 10)$ | .4463 | .4237 | .4636 | .4350 | .4348 | .4326 | .4351 | .4372 | .4182 | .4712 | .4437 | .4712 | .4290 | .4366 | .4443 | .4499 | .4390 | .4345 | .4442 | .4449 |
| 133 $(\times 10)$ | .3099 | .2584 | .3075 | .2980 | .2924 | .2938 | .2964 | .2957 | .2874 | .2971 | .2875 | .3434 | .2943 | .2869 | .2762 | .2739 | .2955 | .3054 | .3455 | .2982 |

| | | | | | | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $135\,(\times 10^3)$ | 4.000 | 5.419 | 9.373 | .6753 | .7938 | 10.16 | 10.15 | 10.12 | 9.828 | 10.35 | 10.01 | 2.725 | 9.682 | 9.970 | .5332 | 1.041 | 8.391 | 4.593 | 5.198 | 3.016 | |
| $138\,(\times 10^{-1})$ | .8865 | 1.205 | .9079 | 1.147 | 1.028 | 1.788 | 1.100 | 1.253 | 2.356 | 1.352 | 1.131 | 3.963 | 2.026 | 1.159 | 1.172 | 1.478 | 2.956 | 2.405 | 3.885 | 1.358 | |
| $140\,(\times 10^{-1})$ | .7825 | .7966 | .7819 | .8101 | .7921 | .7892 | .7908 | .7857 | .7856 | .7883 | .7928 | 1.330 | .7862 | .8258 | .8097 | .8325 | .8596 | .9688 | 1.241 | .7782 | |
| $141\,(\times 10^{-1})$ | .2915 | .2944 | .2925 | .3009 | .2878 | .2973 | .2863 | .2876 | - | .2833 | .2856 | 1.589 | .2942 | .2902 | .3029 | .3195 | .4088 | .4932 | 1.549 | .2848 | |
| $144\,(\times 10^2)$ | .2443 | .2651 | .2241 | .2076 | .2044 | .2383 | .2321 | .2399 | .2287 | .2197 | .2109 | .5529 | .2381 | - | .2043 | .2023 | .3030 | .2264 | .3482 | .2328 | |
| $148\,(\times 10)$ | .3744 | .3874 | 1.627 | .3442 | .3143 | 1.603 | 1.694 | 1.476 | 1.698 | 1.696 | 1.648 | .4047 | 1.651 | 1.700 | .3273 | .3607 | .4628 | .3926 | .6744 | .3946 | |
| 149 | .1390 | .1364 | .1437 | .1325 | .1336 | .1338 | .1369 | .1361 | .1363 | .1360 | .1333 | .3889 | .1355 | .1378 | .1363 | .1337 | .1380 | .1348 | .2294 | .1414 | |
| 152 | .3118 | .3416 | .3554 | .3924 | .3903 | .4795 | .4468 | .4714 | .4417 | .4361 | .3846 | 1.065 | .4332 | - | .4094 | .4249 | .4606 | .4146 | .5252 | .4326 | |
| $153\,(\times 10)$ | .8053 | .9441 | 1.573 | .4826 | .4380 | 1.609 | 1.626 | 1.619 | 1.624 | 1.634 | 1.663 | .6138 | 1.630 | 1.628 | .4658 | .6438 | 1.060 | .9348 | 1.007 | .5151 | |
| $155\,(\times 10)$ | .2416 | .2455 | .2334 | .2404 | .2365 | .2492 | .2331 | .2475 | .2503 | .2441 | .2396 | 1.746 | .2539 | .2463 | .2298 | .2526 | 1.011 | .4781 | 1.711 | .2494 | |
| $156\,(\times 10)$ | .2739 | .2883 | .2779 | .2777 | .2713 | .2911 | .2846 | .2944 | .3058 | .2942 | .2900 | 1.742 | .2941 | .2919 | .2743 | .2934 | 1.021 | .3635 | 1.727 | .2920 | |
| 158 | .3936 | .3940 | .4341 | .3920 | .3895 | .4231 | .4097 | .4150 | .4212 | .4605 | .3981 | .5203 | .4152 | .4026 | .3911 | .3915 | .4851 | .4573 | .5010 | .3962 | |
| $161\,(\times 10)$ | .3971 | .4025 | .4243 | .3967 | .3959 | .3977 | .4081 | .4043 | .4156 | .4334 | .4186 | .6509 | .4012 | .4181 | .4044 | .4070 | .4426 | .4393 | .6487 | .3931 | |
| $162\,(\times 10^{-1})$ | .6253 | .5959 | .7242 | .7196 | .6435 | .6504 | .6356 | .6548 | 1.192 | .6746 | .6881 | 1.492 | .6949 | .6425 | .6409 | .8302 | 1.393 | .6633 | 1.501 | .6776 | |
| $165\,(\times 10^{-2})$ | .1366 | .1372 | .1386 | .1388 | .1383 | .1381 | .1378 | .1389 | .1371 | .1389 | .1376 | .2338 | .1375 | .1374 | .1395 | .1399 | .1423 | .1410 | .2290 | .1379 | |
| $172\,(\times 10^{-2})$ | .1862 | .1966 | .1844 | .2340 | .2030 | .1868 | .1827 | .1850 | .2230 | .1844 | .1868 | .6717 | .1898 | - | .2310 | .2815 | .2899 | .3703 | .6556 | .1867 | |
| 176 | .7874 | .8109 | .7884 | .7721 | .7728 | .8141 | .7978 | .8070 | .7958 | .7991 | .8034 | 1.402 | .8151 | .8033 | .7698 | .7872 | 1.091 | .9369 | 1.312 | .8044 | |
| $178\,(\times 10)$ | .1011 | .1093 | .1017 | .1079 | .1010 | .1022 | .1009 | .1010 | .1007 | .1007 | .1006 | .4949 | .1041 | - | .1063 | .1395 | .2643 | .1854 | .4945 | .1026 | |
| 179 | .1291 | .1487 | .1758 | .1281 | .1315 | .1821 | .1763 | .1765 | .1761 | .1761 | .1762 | .1376 | .1779 | .1763 | .1281 | .1299 | .1570 | .1441 | .1534 | .1542 | |
| 181 | 1.304 | 1.411 | 15.39 | .5704 | .5220 | 15.39 | 15.48 | 16.33 | 15.54 | 15.61 | 15.15 | .6482 | 15.62 | - | .5932 | .7295 | .9492 | 1.014 | 2.165 | .5564 | |
| $184\,(\times 10^4)$ | .4699 | .4492 | .4650 | .4650 | .4736 | .4798 | .4616 | .4874 | .4677 | .5330 | .5255 | 1.310 | .4869 | .4605 | .4696 | .4535 | .6506 | .5315 | 1.284 | .4763 | |
| $188\,(\times 10^5)$ | .3102 | .3452 | .3054 | .3136 | .2994 | .3175 | .3067 | .3133 | .3483 | .3138 | .3216 | .5458 | .3187 | - | .3098 | .3278 | .4819 | .3623 | .5397 | .3158 | |
| $189\,(\times 10^5)$ | .2954 | .3130 | .2922 | .2871 | .2847 | .2971 | .2907 | .2979 | .2996 | .2901 | .2860 | 1.004 | .2981 | .2926 | .2846 | .2891 | .4797 | .3154 | .5345 | .2968 | |
| $190\,(\times 10^6)$ | .1255 | .1490 | .1423 | .1396 | .1177 | .1382 | .1193 | .1309 | .1376 | .1272 | .1368 | .4149 | .1413 | - | .1359 | .1485 | .2406 | .1714 | .3928 | .1394 | |
| $191\,(\times 10^5)$ | .4291 | .4469 | .4214 | .4797 | .4546 | .5200 | .4796 | .5019 | .5426 | .5016 | .4838 | 1.167 | .5205 | .5117 | .4665 | .5247 | .7140 | .6000 | 1.148 | .5128 | |
| $192\,(\times 10^7)$ | .1225 | .1373 | .1108 | .1020 | .1117 | .1113 | .1081 | .1176 | .1045 | .1083 | .1079 | .1510 | .1138 | .1098 | .1027 | .1141 | .1120 | .1121 | .1598 | .1148 | |
| $202\,(\times 10^{-1})$ | .6968 | .8805 | .6639 | 1.321 | .9017 | .6877 | .6736 | .6677 | .9350 | .6559 | .6924 | 4.072 | .7158 | .6857 | 1.215 | 1.492 | 2.054 | 1.205 | 2.620 | .6781 | |
| $205\,(\times 10^{-5})$ | 1830. | 87.82 | 22.92 | 1453. | 1.618 | 751.2 | 24.46 | 63.69 | 63.91 | 92.98 | 81.79 | Inf | 414.3 | 83.63 | .0015 | 0.000 | Inf | Inf | Inf | 205.1 | |
| 212 | .4331 | .4088 | .4440 | .5918 | .4109 | .6367 | .5507 | .5745 | .4221 | .4719 | .5490 | 22.68 | .6340 | .7251 | .4351 | .6579 | 2.353 | 2.487 | 16.72 | .6655 | |
| $223\,(\times 10^{-2})$ | 2.739 | - | 2.797 | 19.28 | 8.548 | 2.867 | 3.095 | 1.624 | .9159 | 2.650 | 3.633 | 1181. | 23.81 | - | 11.38 | 9.087 | 497.5 | 147.8 | 1037. | 9.693 | |
| 233 | .1537 | .1537 | .1341 | .1423 | .1397 | .1393 | .1330 | .1395 | .1321 | - | - | .1329 | .1369 | .1331 | .1412 | .1350 | .1317 | .1416 | .1332 | .1464 | |
| $241\,(\times 10)$ | .2328 | .3402 | .2236 | .4775 | .4189 | .2436 | .2158 | .2208 | .4421 | .2249 | .6091 | 4.051 | .2533 | - | .5038 | .5424 | 3.057 | .8922 | 4.170 | .2520 | |
| $242\,(\times 10)$ | .2185 | .3117 | .2322 | .4383 | .3669 | .2578 | .2216 | .2306 | .4418 | .2365 | .2560 | 5.338 | .2561 | - | .4576 | .4984 | 3.095 | .9333 | 4.170 | .4549 | |
| $244\,(\times 10^{-2})$ | .6083 | .6314 | .5905 | .7873 | .6453 | .8170 | .7535 | .6531 | .8634 | .6826 | .7819 | 2.941 | .8971 | - | .7804 | .7914 | 2.608 | 2.567 | 2.939 | .7009 | |
| $245\,(\times 10)$ | .3258 | .3263 | .3308 | .3304 | .3281 | .3341 | .3285 | .3321 | - | .3306 | .3273 | .5563 | .3366 | .3331 | .3279 | .3272 | .4440 | .3768 | .5517 | .3300 | |
| $247\,(\times 10)$ | .1381 | .1417 | .1172 | .1136 | .1062 | .1162 | .1239 | .1172 | .1122 | .1210 | .1179 | .1731 | .1169 | .1211 | .1244 | .1154 | .1427 | .1112 | .1711 | .1126 | |
| 248 | 1.027 | .9752 | .9483 | .9030 | .8619 | .8886 | .9068 | 1.029 | .8567 | .8900 | .9064 | 1.310 | .8802 | .9175 | .8791 | .8597 | .9821 | .9643 | 1.256 | .9469 | |
| 252 | .6112 | .6935 | .7274 | .7033 | .6682 | .6889 | .7274 | .6992 | .8339 | .6688 | .6973 | 3.361 | .6850 | .6870 | .6900 | .7221 | 1.226 | .7109 | 2.170 | .7165 | |
| 257 | .7146 | .7053 | .7812 | .7245 | .7331 | .7747 | .7234 | .7457 | .7300 | - | .8272 | .8312 | .8148 | .7634 | .7037 | .7073 | .8879 | .7962 | .8328 | .7344 | |
| 258 | .2607 | .2599 | .2974 | .2658 | .2637 | .2658 | .2760 | .2642 | .2622 | .2820 | .2628 | .3244 | .2616 | .2665 | .2639 | .2707 | .3041 | .2641 | .3109 | .2608 | |
| 259 | .3147 | .3168 | .3185 | .3161 | .3127 | .3526 | .3222 | .3452 | .3173 | .3397 | .3197 | .3980 | .3233 | .3329 | .3150 | .3188 | .3882 | .3344 | .4023 | .3197 | |
| $261\,(\times 10^2)$ | .1695 | .1429 | .1982 | .2389 | .2463 | .2802 | .1346 | .2227 | .2754 | .2576 | .1907 | 9.860 | .2097 | .1370 | .2337 | .2501 | .1935 | .2360 | .4149 | .1107 | |
| $262\,(\times 10^{-1})$ | 1.066 | 1.084 | .9945 | 1.138 | 1.091 | .9784 | 1.017 | 1.057 | 1.319 | 1.078 | 1.032 | 2.037 | 1.028 | 1.035 | 1.109 | 1.191 | 1.400 | 1.269 | 1.983 | .9885 | |
| 267 | .6202 | .6084 | .4005 | 1.081 | .7602 | .7415 | .5087 | .5442 | .3696 | .9050 | .5609 | 35.41 | 1.629 | - | 1.067 | 1.353 | 4.958 | 2.945 | 33.25 | .8614 | |
| 269 | .6791 | .6675 | .6899 | .8361 | .7969 | .7572 | .7649 | .8194 | .6973 | .7282 | .7382 | 11.82 | .7656 | .7821 | .9136 | .9651 | 2.132 | .7345 | 6.730 | .9068 | |
| $270\,(\times 10^2)$ | .2059 | .1941 | .1694 | .2236 | .1835 | .1986 | .1887 | .1873 | .1868 | .1779 | .1955 | .1996 | .2058 | .2042 | .1843 | .1802 | .2187 | .1911 | .2000 | .2180 | |
| $271\,(\times 10^{-1})$ | .1986 | .1844 | .1794 | .2185 | .2166 | .2102 | .2127 | .2234 | .3405 | .1659 | .2046 | .5677 | .2164 | .1766 | .2438 | .2471 | .4945 | .2065 | .5669 | .2371 | |
| 272 | .7032 | .7452 | .7410 | .6877 | .6704 | .7739 | .7379 | .7452 | .7066 | .7262 | .7229 | .9334 | .6994 | .7920 | .6829 | .6827 | .8137 | .7717 | .8920 | .7042 | |
| $280\,(\times 10^{-1})$ | .3491 | .3496 | .3579 | .3484 | .3468 | .3505 | .3549 | .6941 | .3497 | .3562 | .3542 | .3578 | .3520 | - | .3484 | .3476 | .3531 | .3522 | .3589 | .3473 | |
| 281 | .2855 | .3670 | .2772 | .2857 | .2816 | .2126 | .3062 | .2185 | .3295 | .2608 | .2692 | 15.71 | .2141 | .2738 | .3103 | .3163 | .2884 | .2989 | 4.244 | .2395 | |
| $286\,(\times 10^2)$ | .4147 | .4877 | .4162 | .3802 | .3821 | .4796 | .4212 | .5103 | .4293 | .3904 | .4253 | .6248 | .4703 | .4547 | .3851 | .3898 | .5630 | .5163 | .6112 | .4748 | |
| $293\,(\times 10)$ | .1202 | .1228 | .1128 | .1254 | .1182 | .1167 | .1170 | .1144 | 1.040 | .1206 | .1173 | 1.438 | .1212 | .1177 | .1329 | .1357 | .1327 | .2006 | 1.303 | .1209 | |
| $295\,(\times 10)$ | .2943 | .2952 | .2913 | .2998 | .2923 | .2884 | .2938 | .2917 | .3477 | .2919 | .2998 | 1.472 | .2933 | .2975 | .3007 | .3084 | .3163 | .3348 | .6349 | .2903 | |
| 297 | .6609 | .6361 | .6895 | .6613 | .6569 | .7314 | .7287 | .7441 | - | .7061 | .7240 | .9099 | .7470 | .7361 | .6751 | .6758 | .7690 | .6793 | .9149 | .7192 | |