



**ICLR**

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

**Han-Jia Ye, Huai-Hong Yin, De-Chuan Zhan**  
Nanjing University

**Wei-Lun (Harry) Chao**  
The Ohio State University

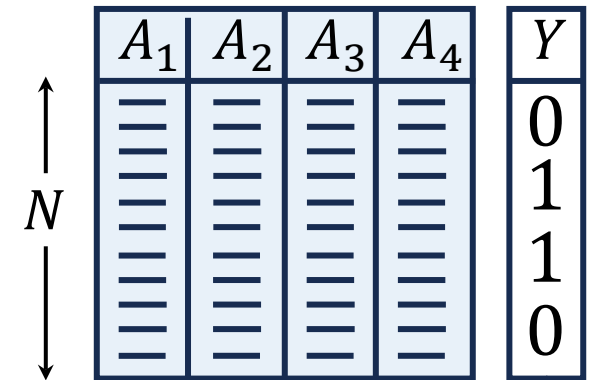
Paper link: <https://openreview.net/forum?id=JytL2MrILT>  
Code : <https://github.com/LAMDA-Tabular/TALENT>

# Learning with Tabular Data

- Given a training set  $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$ ,
  - $N$  instances (rows) and  $d$  columns (attributes)
  - $\mathbf{x}_i \in \mathbb{R}^d$ , with *categorical* and *numerical* features/attributes
  - $y_i \in \{0,1\}$  for binary classification,  $y_i \in \{1, \dots, C\}$  for C-way classification, and  $y_i \in \mathbb{R}$  for regression
- The goal is to learn a mapping  $f$ . Given an unseen instance  $\mathbf{x}^*$ ,

$$\hat{y}^* = f(\mathbf{x}^*, \mathcal{D})$$

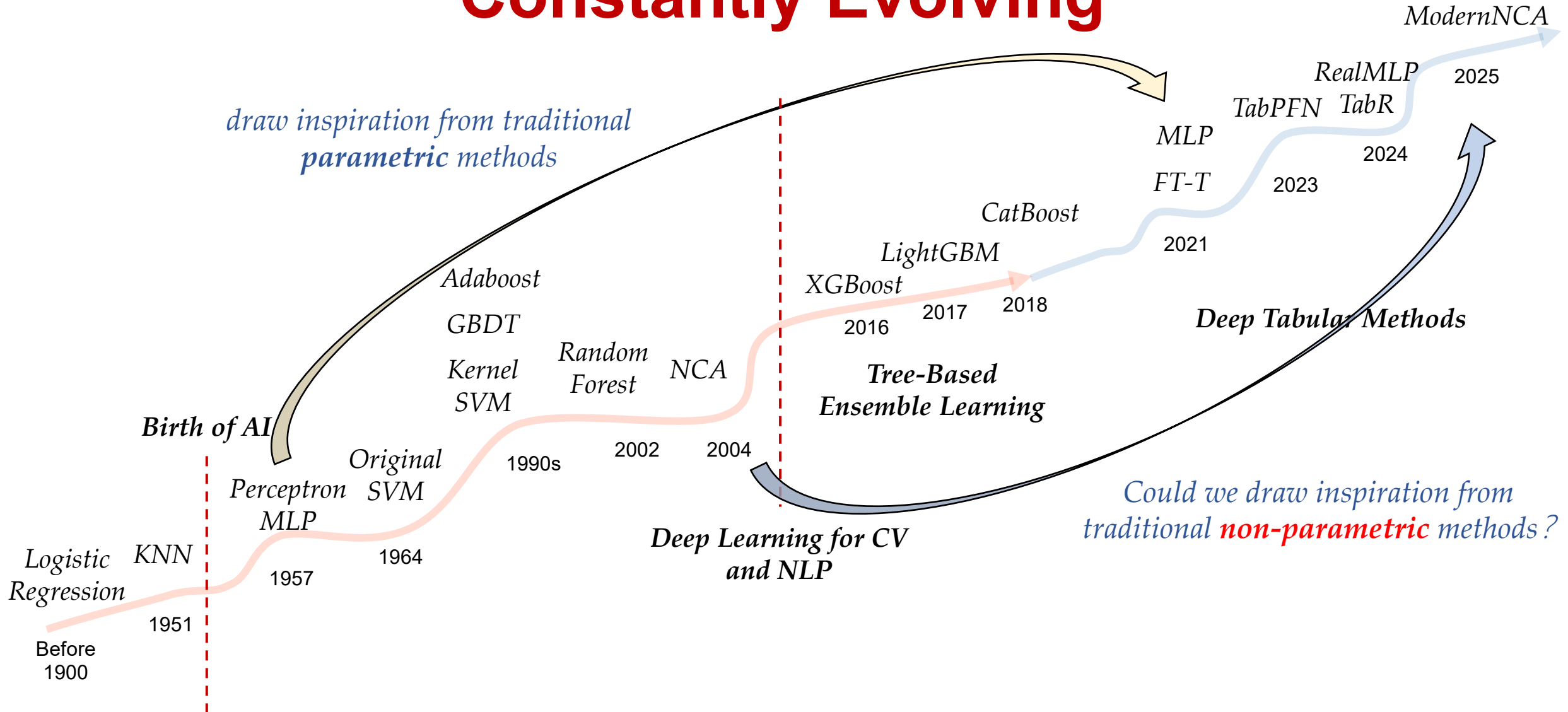
- Other related tabular tasks
  - Clustering, imputation, generation, anomaly detection, etc.



$A_1$	$A_2$	$A_3$	$A_4$	$Y$
—	—	—	—	0
—	—	—	—	1
—	—	—	—	1
—	—	—	—	0

*Binary Classification*

# Methods for Tabular Datasets are Constantly Evolving



# Neighborhood Component Analysis (NCA)

- Neighborhood Component Analysis (NCA) for classification [Goldberger et al., NIPS'04]
- The probability that  $\mathbf{x}_j$  locates in the neighborhood of  $\mathbf{x}_i$

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

- The probability that an instance  $\mathbf{x}_i$  is classified as the class  $y_i$  is

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

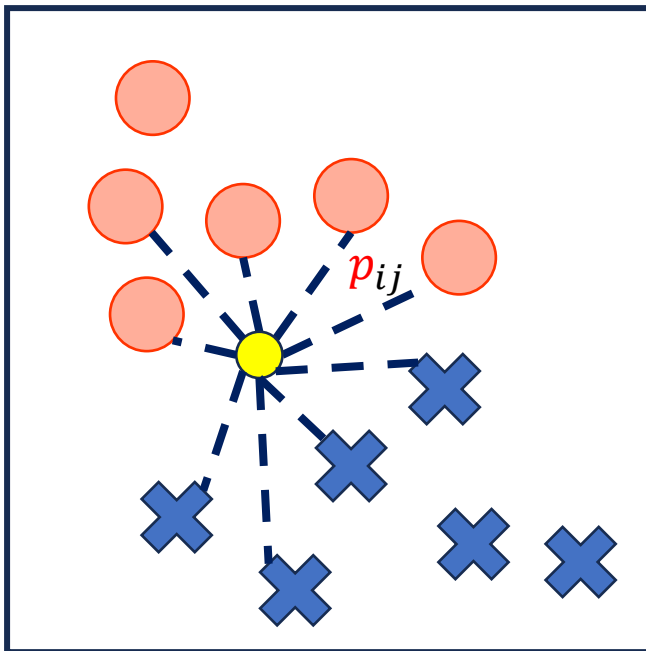
- NCA maximizes the sum of  $\Pr(\hat{y}_i = y_i \mid \mathbf{x}_i, \mathcal{D}, \mathbf{L})$  over all instance in  $\mathcal{D}$  and use KNN over the space projected by  $\mathbf{L}$  in the test phase.

# Our First Attempt

- The learning objective: soft-NN

$$\hat{y}_i = \sum_{(x_j, y_j) \in \mathcal{D}} p_{ij} y_j = \sum_{(x_j, y_j) \in \mathcal{D}} \frac{\exp\left(-\text{dist}^2\left(\phi(x_i), \phi(x_j)\right)\right)}{\sum_{(x_l, y_l) \in \mathcal{D}, x_l \neq x_i} \exp\left(-\text{dist}^2\left(\phi(x_i), \phi(x_l)\right)\right)} y_j$$

- $\phi(x_i) = L^\top x_i$ . Minimize the sum of negative **log probability**.



- Prediction Strategy

- Use Soft-NN to make prediction
- Do not restrict  $L$  to project into low-dimensional space
- Use SGD instead of L-BFGS

Denote this improved *linear* version of NCA as **L-NCA**

# Our Second Attempt

- Architectures

- Define  $\phi$  as MLP, and a single block is implemented by [Gorishniy et al., NIPS'21]

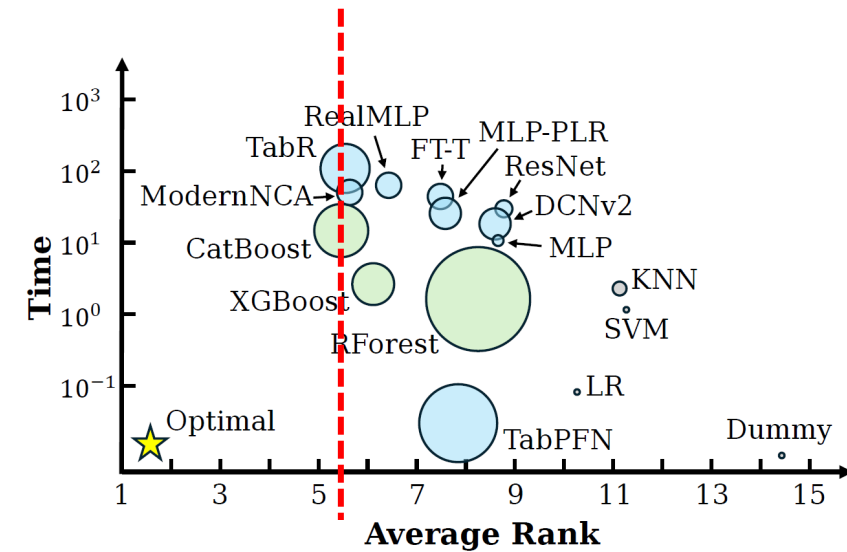
$$g(x_i) = \text{Linear}(\text{Dropout}((\text{ReLU}(\text{Linear}(\text{BatchNorm}(x_i))))))$$

- One-hot encoding for categorical features
  - PLR (lite) encoding for numerical features
- Stochastic Neighborhood Sampling (SNS)
  - Sample *a subset of training set* in a mini-batch to act as the neighbor candidates, while use *the whole training set* during the inference.
- Distance Function: Euclidean distance

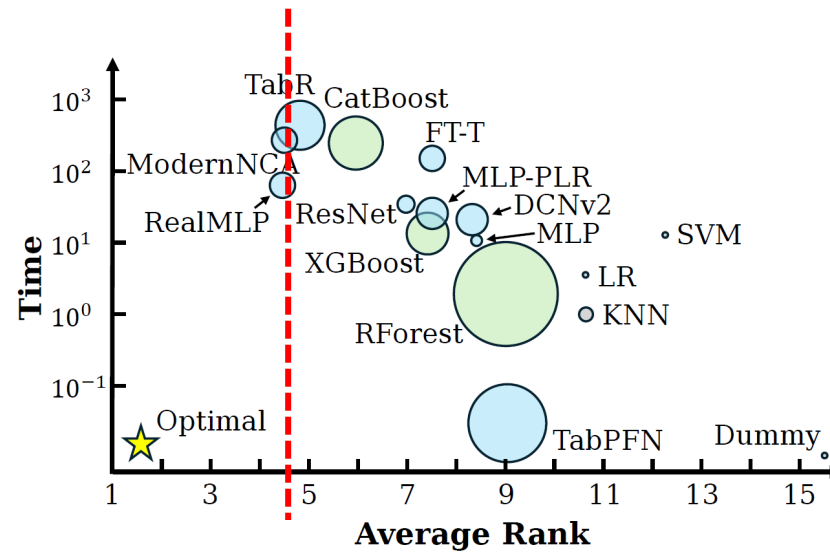
Denote this improved *nonlinear* version of NCA as **ModernNCA (MNCA)**

# Experiments

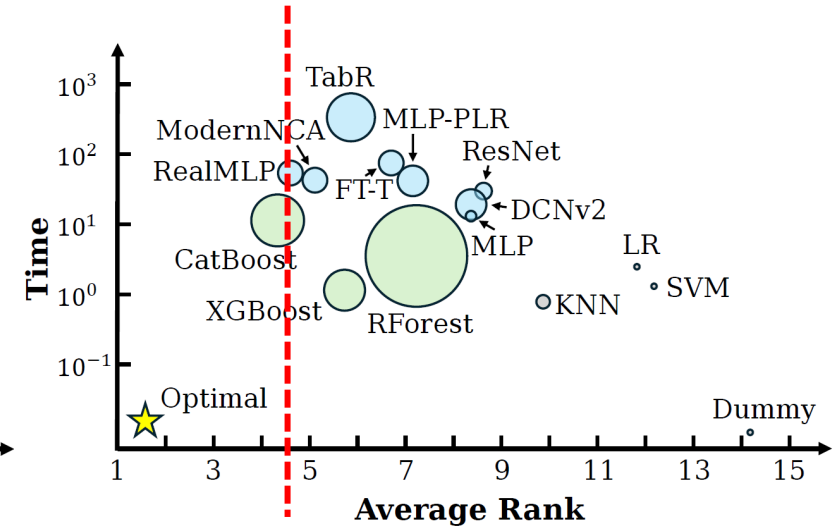
- 300 datasets in total
  - 120 binary classification, 80 multi-class classification, 100 regression
- Average rank vs. Training time vs. Model Size



*Binary Classification*



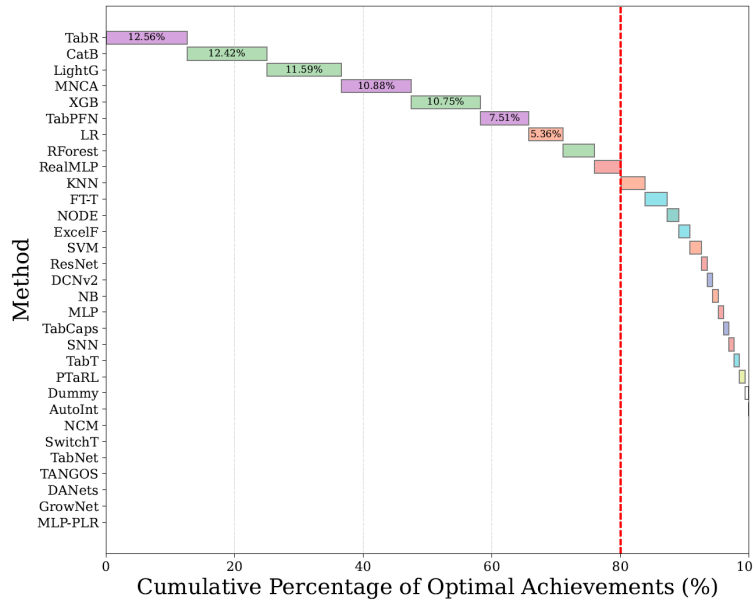
*Multi-Class Classification*



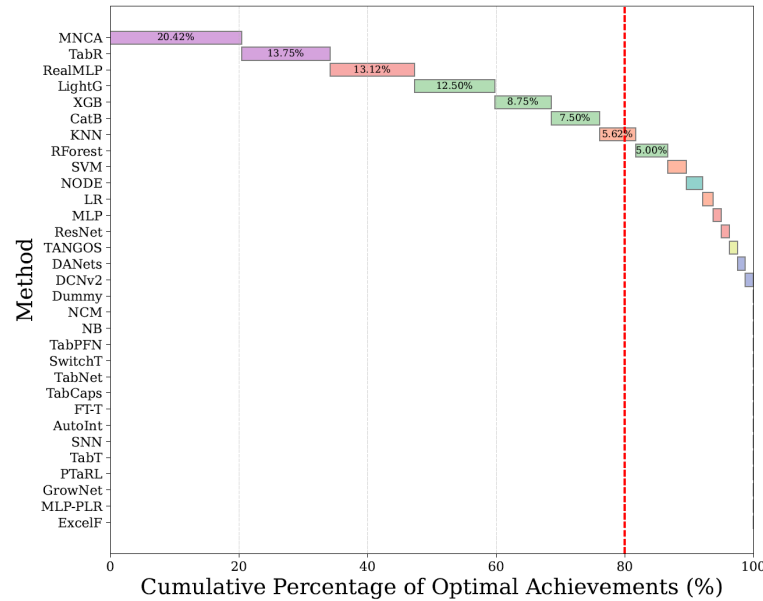
*Regression*

# Experiments

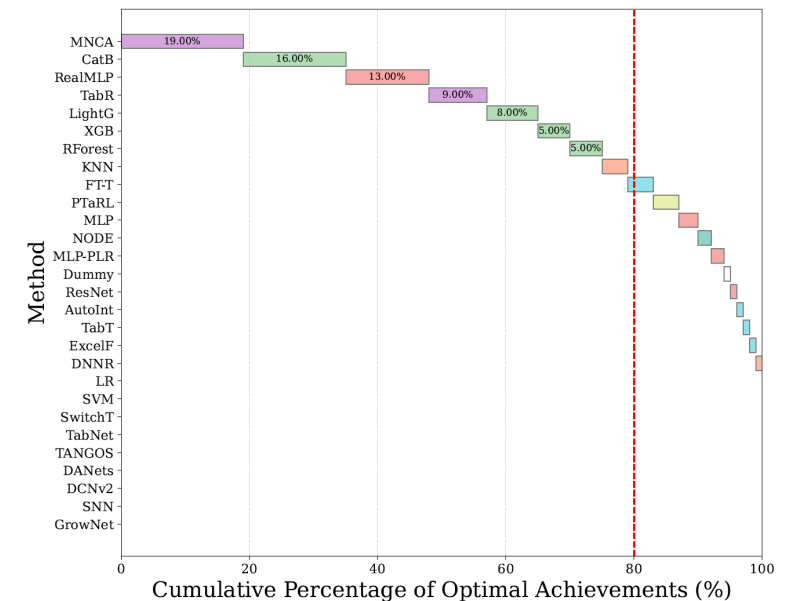
- 300 datasets in total
  - 120 binary classification, 80 multi-class classification, 100 regression
- PAMA (Probability of Achieving the Best Accuracy)



*Binary Classification*



*Multi-Class Classification*



*Regression*



# Ablations

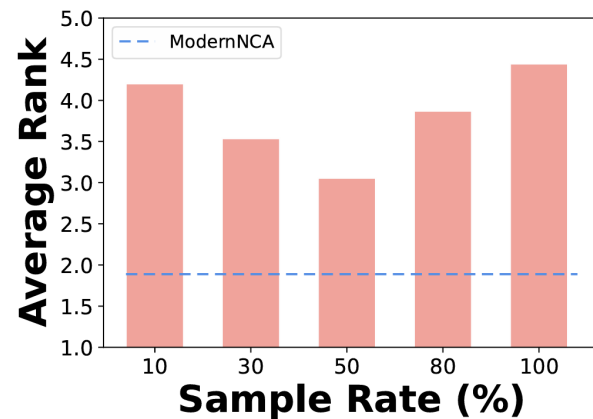
- Tiny-Benchmark with 45 datasets
- From NCA to L-NCA on 27 classification datasets

	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

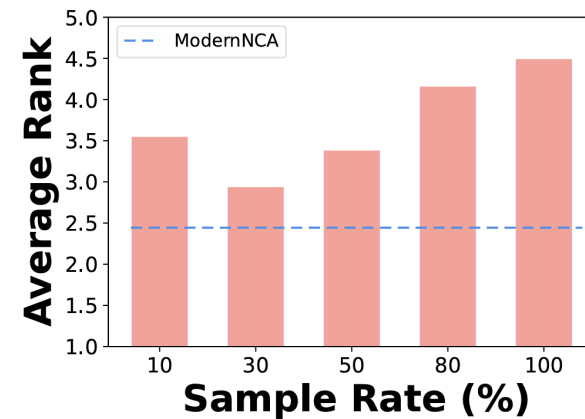
# Ablations

- Tiny-Benchmark with 45 datasets
- Average rank changes from L-NCA to M-NCA

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



(a) classification



(b) regression

# Come to out poster!

## Poster Session 3

Fri 25 Apr 10 a.m. CST – 12:30 p.m. CST

- We revisit and enhance one of the most representative neighborhood-based methods, NCA, by incorporating modern deep learning techniques.
- **ModernNCA** establishes itself as a very strong baseline for deep tabular prediction, frequently outperforming both tree-based and deep learning models across a wide range of classification and regression tasks.

*For more discussions, please contact [yehj@nju.edu.cn](mailto:yehj@nju.edu.cn)*

# Tabular Toolbox



TALENT: A Tabular Analytics and Learning Toolbox

<https://github.com/LAMDA-Tabular/TALENT>

30+ deep learning methods (with methods on NeurIPS'24/ICLR'25), unifying interfaces, customizability.