DP-KAN: Differentially Private Kolmogorov-Arnold Networks

Anonymous Full Paper Submission 13

OD1 Abstract

We study the Kolmogorov-Arnold Network (KAN), 002 recently proposed as an alternative to the classical 003 Multilayer Perceptron (MLP), in the application for 004 005 differentially private model training. Using the DP-SGD algorithm, we demonstrate that KAN can be 006 made private in a straightforward manner and eval-007 uated its performance across several datasets. Our 008 results indicate that the accuracy of KAN is not only 009 comparable with MLP but also experiences similar 010 deterioration due to privacy constraints, making it 011 suitable for differentially private model training. 012

1 Introduction

The Kolmogorov-Arnold Network (KAN) [1] has 014 recently emerged as a new approach to symbolic 015 regression and general prediction problems. This 016 new architecture already saw remarkable attention 017 in very recent papers [2-6]. Given its promising 018 capabilities, we find it important to test the ability 019 of KAN to be trained in a privacy-preserving man-020 ner without compromising sensitive information. In 021 this work, we chose differential privacy as a way to 023 protect privacy.

As data analysis, and particularly machine learn-024 ing, continues to expand, the demand for more data 025 increases. However, there is a growing concern about 026 the potential misuse or exploitation of personal in-027 formation. Differential Privacy [7] is a time-tested 028 notion of formally defined privacy of algorithms 029 concerning sensitive data. Over the past nearly two 030 decades, it has attracted considerable attention from 031 both theoretical and practical perspectives. 032

The primary result of this work is that we pro-033 vide the first integration of KANs with dif-034 ferentially private training algorithms, specif-035 ically with DP-SGD. Our study shows that non-036 private KANs and MLPs exhibit comparable accu-037 racy. Moreover, when trained with differential pri-038 vacy, KANs experience similar accuracy degradation 039 compared to MLPs. This makes KANs a promis-040 ing option for maintaining model performance while 041 ensuring differential privacy in training. 042

043 1.1 Related Works

Kolmogorov Arnold Networks (KANs) have recently
gained significant attention, with various studies
building upon the original framework introduced by

Ziming et al. [1]. Zavareh and Chen [6] introduced 047 Wav-KAN, enhancing interpretability and perfor-048 mance by integrating wavelet functions to efficiently 049 capture both high-frequency and low-frequency com-050 ponents of input data. For time series forecasting, 051 Genet and Inzirillo [2] proposed Temporal KAN 052 (TKAN), combining the strengths of Long Short-053 Term Memory (LSTM) networks and KANs, while 054 Vaca-Rubio et al. [4] demonstrated that KANs out-055 perform conventional MLPs in satellite traffic fore-056 casting with fewer parameters. Xu et al. [5] intro-057 duced FourierKAN-GCF, a graph-based recommen-058 dation model that utilizes a Fourier KAN to improve 059 feature transformation during message passing in 060 graph convolution networks (GCNs), achieving su-061 perior performance in recommendation tasks. A 062 benchmarking study by Poeta et al. [8] on real-063 world tabular datasets shows that KANs achieve 064 superior or comparable accuracy and F1 scores com-065 pared to MLPs, especially on larger datasets, though 066 with higher computational costs. Other notable ad-067 vancements include Shukla et al. [3], which com-068 pared KANs with traditional Multi-Layer Percep-069 trons (MLPs), and Basim and Naveed [9], which 070 developed Convolutional KAN (ConvKAN) for en-071 hanced image processing. 072

Differentially private stochastic gradient descent 073 (DP-SGD) has become a standard in private ma-074 chine learning, being applied in both convex and 075 non-convex optimization. Bassily et al. [10] demon-076 strated its effectiveness in convex optimization, while 077 Abadi et al. [11] extended its application to deep 078 learning, ensuring privacy in training deep neural 079 networks. In our study, we use DP-Adam, a specific 080 member of the DP-SGD family of algorithms, to 081 leverage its adaptive learning rate capabilities for 082 improved performance. The DP-Adam method has 083 been studied in various works, including Gylberth 084 et al. [12], who demonstrated improved accuracy 085 and faster convergence, and Tang et al. [13], who 086 proposed DP-AdamBC to correct bias in the second 087 moment estimation. In differentially private regres-088 sion, our tabular data experiments build on the work 089 by Amin et al. [14], who introduced a method using 090 the exponential mechanism to select a model with 091 high Tukey depth, eliminating the need for data 092 bounds or hyperparameter selection. Other signif-093 icant contribution includes Alabi et al. [15], who 094 designed differentially private algorithms for simple 095 linear regression in small datasets. 096

⁰⁹⁷ 2 Background

2.1 Kolmogorov-Arnold Networks

Kolmogorov-Arnold Networks (KAN) [1] are a novel neural network architecture inspired by the Kolmogorov-Arnold representation theorem [16], which states that any multivariate continuous function $f: [0,1]^n \to \mathbb{R}$ can be decomposed into a sum of compositions of univariate functions:

105
$$f(x^{(1)}, \dots, x^{(n)}) = \sum_{q=1}^{2n+1} \Phi_q \left(\sum_{p=1}^n \phi_{q,p}(x^{(p)}) \right).$$
 (1)

KANs replace fixed activation functions with
learnable univariate functions parameterized as Bsplines. This enhances flexibility and interpretability compared to traditional Multi-Layer Perceptrons
(MLPs), where activations occur at the nodes and
weights are linear.

112 In KANs, B-splines serve as the building blocks 113 for learnable univariate functions $\phi_{q,p}$, expressed as:

114
$$\operatorname{spline}(x) = \sum_{i} c_i B_i(x),$$
 (2)

where c_i are coefficients learned during training and $B_i(x)$ are B-spline basis functions. KANs use residual activation functions, combining a basis function b(x) with the B-spline:

$$\phi(x) = w_b b(x) + w_s \text{spline}(x), \qquad (3)$$

where b(x) is a sigmoid linear unit (SiLU), and w_b and w_s are trainable weights.

KANs are trained using backpropagation, compatible with standard optimization techniques such as Differentially Private SGD (DP-SGD 1). They have shown superior performance to MLPs in specific tasks, such as formula recovery and symbolic regression [1].

128 2.2 Differential Privacy

129 Differential privacy guarantees that a (random-130 ized) algorithm is *stable* with respect to single data 131 point changes of the input dataset. In particular, 132 we say that two datasets D and D' are adjacent if 133 they differ only in one sample. Then, we have the 134 following

135 Definition 2.1 ((ε, δ) -differential privacy [7]) **136** A randomized algorithm $\mathcal{A} : \mathcal{D} \to \mathcal{S}$ satisfies **137** (ε, δ)-differential privacy if for any two adjacent **138** inputs $D, D' \in \mathcal{D}$, and for any subset of the output **139** space $S \subseteq \mathcal{S}$, we have

140
$$\mathbb{P}\left(\mathcal{A}(D)\in S\right) \le e^{\varepsilon}\mathbb{P}\left(\mathcal{A}(D')\in S\right) + \delta.$$
(4)

Here, \mathcal{D} represents the space of all datasets, and \mathcal{S} , in the case of deep learning models, is the space of the

Algorithm 1 Differentially private gradient descent with Adam Optimizer

- **Input:** Number of iterations T, learning rate η , clipping constant C, noise multiplier σ , batch size B, initialization θ_0 .
 - 1: for $t \in [T]$ do
 - 2: Randomly sample a batch of B samples
- 3: Compute the gradients $g(x_i, y_i, \theta_{t-1}) \leftarrow \nabla_{\theta} \ell(x_i, y_i, \theta_{t-1})$

4: Clip the gradients

$$\tilde{g}(x_i, y_i, \theta_{t-1}) \leftarrow g(x_i, y_i, \theta_{t-1})/\zeta,$$

where $\zeta = \max\left(1, \frac{\|g(x_i, y_i, \theta_{t-1})\|_2}{C}\right)$

5: Aggregate the noisy gradients

$$g_t \leftarrow \frac{1}{B} \sum_i \tilde{g}(x_i, y_i, \theta_{t-1}) + \frac{\sigma C}{B} \mathcal{N}(0, I)$$

 $\theta_t = \operatorname{Adam}(\theta_{t-1}, g_t)$

7: end for

6:

Output: Model parameters θ_T

trainable parameters. Differential privacy ensures 143 that from the resulting model, one cannot extract 144 the training data with high probability, nor can one 145 determine with confidence which specific data points 146 were used during training. This property is crucial 147 for protecting individual privacy in machine learning 148 applications, particularly in contexts where sensitive 149 data may be involved. 150

In machine learning, a well-established family of 151 algorithms allows us to train models with differen-152 tial privacy (DP) guarantees: differentially private 153 stochastic gradient descent (DP-SGD). These algo-154 rithms [11] involve an iterative process where gradi-155 ents are clipped to have a bounded maximum norm 156 and then summed up with Gaussian noise scaled by 157 a *noise multiplier* to ensure given privacy guaran-158 tees. The algorithm described in 1 builds on the 159 Adam optimizer, illustrating one specific example 160 within this family. 161

3 Results

162

We conducted two sets of experiments: regression on 163 tabular data and classification on the MNIST and 164 USPS image datasets. For the regression task, we 165 used various tabular datasets 1 from Amin et al. [14] 166 and trained differentially private and non-private 167 models. Specifically, we used mean squared error 168 and a one-layer neural network model for linear re-169 gression and one layer KAN in our experiments. We 170 used the datasets and training settings from Amin 171 et al. [14]. On each of the datasets and each of 172 the models, we computed the coefficient of deter-173 mination $(R^2 \text{ score})$. The results can be found in 174 Table 1 and the hyperparameters in Table A.3. KAN 175

 $^{^1\}mathrm{We}$ were unable to reproduce the results for the Beijing dataset, so it has not been included in our analysis.

216



Figure 1. Validation accuracy vs the number of parameters for FasterKAN, MLP, DP MLP, and DP FasterKAN on MNIST and USPS datasets. The privacy garantees for the MNIST models are $(0.87, 10^{-5})$ DP and $(2.03, 10^{-5})$ DP for USPS dataset. The error bars are based on three trials for each point.

model demonstrated lower quality degradation dueto privacy across most datasets.

We also trained differentially private and non-178 private models on the MNIST and USPS datasets 179 using the fasterKAN [17]. This implementation has 180 a superior computational speed compared to the 181 official pykan implementation [1], which was found 182 to be inefficient and impractical for large datasets. 183 For differential privacy, we employed the DP Adam 184 optimizer via the pyvacy library [18], an unofficial 185 PyTorch adaptation of TensorFlow Privacy, which 186 was more PyTorch compatible with fasterKAN [17]. 187 We compared the quality of fasterKAN against a 188 Multi-Layer Perceptron (MLP) in a setting of two 189 layers networks with varying hidden layer sizes, re-190 191 sulting in different numbers of parameters. We used CrossEntropyLoss as the loss function and the ac-192 curacy on the test dataset for evaluation, the hy-193 perparameters of those experiments can be found in 194 the Table A.4. The results, shown in Figure A.1, in-195 dicate that fasterKAN consistently achieved higher 196 accuracy for a relatively lower number of parameters 197 compared to the MLP models. In the differentially 198 private setting, fasterKAN suffered a similar quality 199 degradation to the MLP models. This demonstrates 200 the effectiveness of fasterKAN in balancing accuracy 201 and privacy constraints, making it an appropriate 202 choice for differentially private training scenarios. 203

²⁰⁴ 4 Conclusion

In this study, we showed that the KolmogorovArnold Network (KAN) is a possible alternative to
the classical Multi-Layer Perceptron (MLP) for differentially private training scenarios. Through experiments on both tabular data and MNIST and USPS

image classification tasks, KAN not only achieves 210 comparable accuracy with MLP but also shows similar performance under privacy constraints. Future 212 work can explore further optimizations and applications of KAN in various privacy-preserving machine 214 learning contexts. 215

References

- L. Ziming, Y. Wang, S. Vaidya, F. Ruehle, 217
 J. Halverson, M. Soljačić, T. Hou, and M. 218
 Tegmark. KAN: Kolmogorov-Arnold networks. 219
 arXiv preprint arXiv:2404.19756. 2024. DOI: 220
 10.32388/7NNCAA. 221
- R. Genet and H. Inzirillo. TKAN: Temporal 222 Kolmogorov-Arnold Networks. arXiv preprint 223 arXiv:2405.07344. 2024. DOI: 10.48550 / 224 arXiv.2405.07344. 225
- [3] K. Shukla, J. D. Toscano, Z. Wang, Z. Zou, and 226
 G. E. Karniadakis. A comprehensive and FAIR 227
 comparison between MLP and KAN represen-228
 tations for differential equations and opera-229
 tor networks. arXiv preprint arXiv:2406.02917. 230
 2024. DOI: 10.2139/ssrn.4858126. 231
- [4] C. J. Vaca-Rubio, L. Blanco, P. R., and M. 232 Caus. Kolmogorov-Arnold networks (KANs) 233 for time series analysis. arXiv preprint 234 arXiv:2405.08790. 2024. DOI: 10.48550 / 235 arXiv.2405.08790. 236
- J. Xu, Z. Chen, J. Li, S. Yang, W. Wang, 237
 X. Hu, and E. C. H. Ngai. FourierKAN- 238
 GCF: Fourier Kolmogorov-Arnold Network- 239
 An Effective and Efficient Feature Transforma- 240
 tion for Graph Collaborative Filtering. arXiv 241

Dataset	Linear Reg	KAN	DP-SGD Reg	DP-KAN	Lin Reg drop, $\%$	KAN drop, $\%$
Synthetic	0.997	0.996	0.997 ± 0.001	0.953 ± 0.035	0.1	4.3
California	0.637	0.652	0.091 ± 0.009	0.172 ± 0.020	85.7	73.7
Diamonds	0.907	0.926	0.829 ± 0.000	0.924 ± 0.000	8.6	0.3
Traffic	0.966	0.979	0.937 ± 0.002	0.940 ± 0.000	3.0	4.0
NBA	0.622	0.631	0.529 ± 0.051	0.572 ± 0.015	15.0	9.26
Garbage	0.546	0.554	0.275 ± 0.027	0.371 ± 0.148	49.6	33.0
MLB	0.722	0.723	0.714 ± 0.004	0.719 ± 0.002	1.1	0.6

Table 1. R^2 scores for different models in the experimental setup of Amin et al. [14]. For each of the datasets we computed the relative drop in quality due to privacy. We provide error intervals based on three trials, calculated as half the difference between the maximum and minimum values. We also observe that there is almost no stochasticity in non-private models; therefore, we do not include errors for them.

- 242
 preprint arXiv:2406.01034.
 2024.
 DOI:
 10.
 [14]
 48550/arXiv.2406.01034.
- 244 [6] B. Zavareh and H. Chen. Wav-KAN: Wavelet
 245 Kolmogorov-Arnold Networks. arXiv preprint
 246 arXiv:2405.12832. 2024. DOI: 10.48550 /
 247 arXiv.2405.12832.
- [7] C. Dwork. "Differential privacy". In: International Colloquium on Automata, Languages, and Programming (ICALP). 2006. DOI: 10.
 1007/11787006_1.
- [8] E. Poeta, F. Giobergia, E. Pastor, T. Cerquitelli, and E. Baralis. A Benchmarking
 Study of Kolmogorov-Arnold Networks on Tabular Data. arXiv preprint arXiv:2406.14529.
 2024. DOI: 10.48550/arXiv.2406.14529.
- [9] A. Basim and A. Naveed. Suitability of KANs for Computer Vision: A preliminary investigation. arXiv preprint arXiv:2406.09087. 2024.
 DOI: 10.48550/arXiv.2406.09087.
- [10] R. Bassily, A. Smith, and A. Thakurta. "Private empirical risk minimization: Efficient algorithms and tight error bounds". In: *IEEE Symposium on Foundations of Computer Science*(FOCS). 2014. DOI: 10.1109/F0CS.2014.56.
- [11] M. Abadi, A. Chu, I. Goodfellow, H. B. McMahan, I. Mironov, K. Talwar, and L. Zhang.
 "Deep learning with differential privacy". In: *ACM Conference on Computer and Communications Security (CCS)*. 2016. DOI: 10.1145/
 2976749.2978318.
- [12] R. Gylberth, R. Adnan, S. Yazid, and T. Basaruddin. "Differentially private optimization algorithms for deep neural networks". In: *Advanced Computer Science and Information Systems (ICACSIS)* (2017). DOI: 10.1109/
 ICACSIS.2017.8355063.
- 278 [13] Q. Tang, F. Shpilevskiy, and M. Lécuyer.
 279 "Your DP-Adam Is Actually DP-SGD (Unless
 280 You Apply Bias Correction)". In: Conference
 281 on Artificial Intelligence (AAAI) (2024). DOI:
 282 10.1609/aaai.v38i14.29451.

- K. Amin, M. Joseph, M. Ribero, and S. Vassilvitskii. "Easy Differentially Private Linear 284 Regression". In: International Conference on 285 Learning Representations (ICLR). 2022. DOI: 286 10.2478/popets-2022-0041. 287
- D. Alabi, A. McMillan, J. Sarathy, A. Smith, 288
 and S. Vadhan. "Differentially Private Simple 289
 Linear Regression." In: *Privacy Enhancing 290 Technologies Symposium (PETS)*. 2022. DOI: 291
 10.2478/popets-2022-0041. 292
- [16] A. Kolmogorov. "On the representation of continuous functions of several variables by superpositions of continuous functions of a smaller number of variables". In: *American Mathematical Society* (1961). DOI: 10.1007/978-3-642-01742-1_5.
- [17] A. Delis. FasterKAN. https://github.com/ 299 AthanasiosDelis/faster-kan/. 2024. 300
- [18] C. Waites. pyvacy. https://github.com/ 301 ChrisWaites/pyvacy. 2019. 302

A Appendix

303

#13



Figure A.1. Validation accuracy versus the number of parameters on MNIST for FasterKAN, MLP, DP MLP, and DP FasterKAN at different noise levels $\sigma \in \{0.5, 0.25\}$, which corresponds to $(\{7.5, 121\}, 10^{-5})$ Differential Privacy.

Model	Width	Parameters	Test Accuracy	DP Test Accuracy			
				$\sigma = 1$	$\sigma = 0.5$	$\sigma = 0.25$	
	2048	3,257,888	0.970 ± 0.008	0.754 ± 0.018	0.814 ± 0.005	0.840 ± 0.003	
	1024	1,629,728	0.970 ± 0.004	0.775 ± 0.007	0.826 ± 0.004	0.869 ± 0.002	
	512	815,648	0.966 ± 0.005	0.775 ± 0.012	0.832 ± 0.001	0.876 ± 0.002	
KAN	256	408,608	0.967 ± 0.001	0.775 ± 0.017	0.833 ± 0.006	0.884 ± 0.007	
MAN	128	205,088	0.967 ± 0.004	0.777 ± 0.020	0.832 ± 0.005	0.880 ± 0.000	
	64	103, 328	0.958 ± 0.007	0.774 ± 0.007	0.829 ± 0.010	0.882 ± 0.006	
	32	52,448	0.943 ± 0.002	0.783 ± 0.005	0.828 ± 0.012	0.881 ± 0.002	
	16	27,008	0.907 ± 0.012	0.757 ± 0.007	0.827 ± 0.008	0.872 ± 0.008	
	4096	3,256,330	0.968 ± 0.004	0.798 ± 0.006	0.832 ± 0.005	0.873 ± 0.005	
	2048	1,628,170	0.967 ± 0.007	0.798 ± 0.013	0.798 ± 0.012	0.829 ± 0.023	
	1024	814,090	0.965 ± 0.005	0.791 ± 0.010	0.817 ± 0.005	0.846 ± 0.010	
	512	407,050	0.961 ± 0.004	0.801 ± 0.016	0.837 ± 0.005	0.882 ± 0.008	
MLP	256	203, 530	0.960 ± 0.002	0.785 ± 0.018	0.837 ± 0.011	0.896 ± 0.002	
	128	101,770	0.948 ± 0.004	0.784 ± 0.012	0.848 ± 0.006	0.906 ± 0.003	
	64	50,890	0.930 ± 0.004	0.790 ± 0.006	0.840 ± 0.002	0.907 ± 0.005	
	32	25,450	0.881 ± 0.009	0.783 ± 0.005	0.833 ± 0.010	0.900 ± 0.005	
	16	12,730	0.810 ± 0.032	0.767 ± 0.007	0.810 ± 0.003	0.899 ± 0.001	

Table A.1. Test Accuracy for FasterKAN, MLP, DP MLP, and DP FasterKAN on MNIST based on 3 trials.

Model	Width	Parameters	Test Accuracy	DP Test Accuracy	
	16	4,282	0.837 ± 0.019	0.661 ± 0.023	
	32	8,554	0.860 ± 0.005	0.689 ± 0.009	
	64	17,098	0.877 ± 0.005	0.703 ± 0.001	
	128	34, 186	0.897 ± 0.006	0.702 ± 0.004	
MLP	256	68,362	0.909 ± 0.007	0.715 ± 0.008	
	512	136,714	0.920 ± 0.005	0.721 ± 0.017	
	1024	273,418	0.928 ± 0.003	0.740 ± 0.005	
	2048	546,826	0.928 ± 0.008	0.747 ± 0.011	
	4096	1,093,642	0.934 ± 0.006	0.766 ± 0.005	
	16	9,056	0.894 ± 0.005	0.684 ± 0.012	
	32	17,600	0.917 ± 0.007	0.699 ± 0.022	
KAN	64	34,688	0.929 ± 0.005	0.729 ± 0.013	
	128	68,864	0.944 ± 0.004	0.737 ± 0.008	
	256	137,216	0.942 ± 0.001	0.747 ± 0.011	
	512	273,920	0.949 ± 0.004	0.757 ± 0.011	
	1024	547,328	0.945 ± 0.002	0.764 ± 0.011	
	2048	1,094,144	0.936 ± 0.013	0.757 ± 0.004	

Table A.2. Test Accuracy for FasterKAN, MLP, DP MLP, and DP FasterKAN on USPS dataset based on 3 trials.

Dataset	Noise mult.	DP KAN			DP Linear Reg				
		Epochs	C. N.	Lr	Bs	Epochs	C. N.	Lr	Bs
Synthetic	1.472	20	1	1	128	20	1	1	128
California	1.178	20	100	1	64	20	100	1	64
Diamonds	1.089	20	10^{6}	1	128	20	10^{6}	1	128
Traffic	2.016	1	10^{6}	1	1024	1	10^{6}	1	1024
NBA	2.468	20	100	1	512	20	100	1	512
Garbage	0.998	20	1	1	32	20	1	1	32
MLB	1.066	10	100	0.01	512	10	100	0.01	512

Table A.3. The table presents the hyperparameters for tabular data experiments, including the Noise multiplier, number of epochs, the Clipping Norm (C. N.), the learning rate (Lr), and the batch size (Bs), for both Adam and DP Adam optimizers. For KAN, we use a grid size of 2 and a value of k equal to 2. We ensure $(\log(3), 10^{-5})$ -differential privacy in the same manner as described in [14].

Method	Epochs	Batch Size	Noise Multiplier	Clipping Norm	Learning Rate
FasterKAN	15	64	0	-	0.001
MLP	15	64	0	-	0.001
DP FasterKAN	15	64	$\{1, 0.5, 0.25\}$	10^{-3}	0.001
DP MLP	15	64	$\{1, 0.5, 0.25\}$	10^{-3}	0.001

Table A.4. Hyperparameters for MNIST and USPS experiments. Number of epochs, batch size, noise multiplier, max grad norm, learning rate. Trained with AdamW and DP Adam. FasterKAN with parameters grid_min =-1.2, grid_max = 0.2, num_grids = 2, exponent = 2, inv_denominator = 0.5, train_grid =False, train_inv_denominator=False. We used batch clipping for non-private training to ensure the stability of the optimization.