

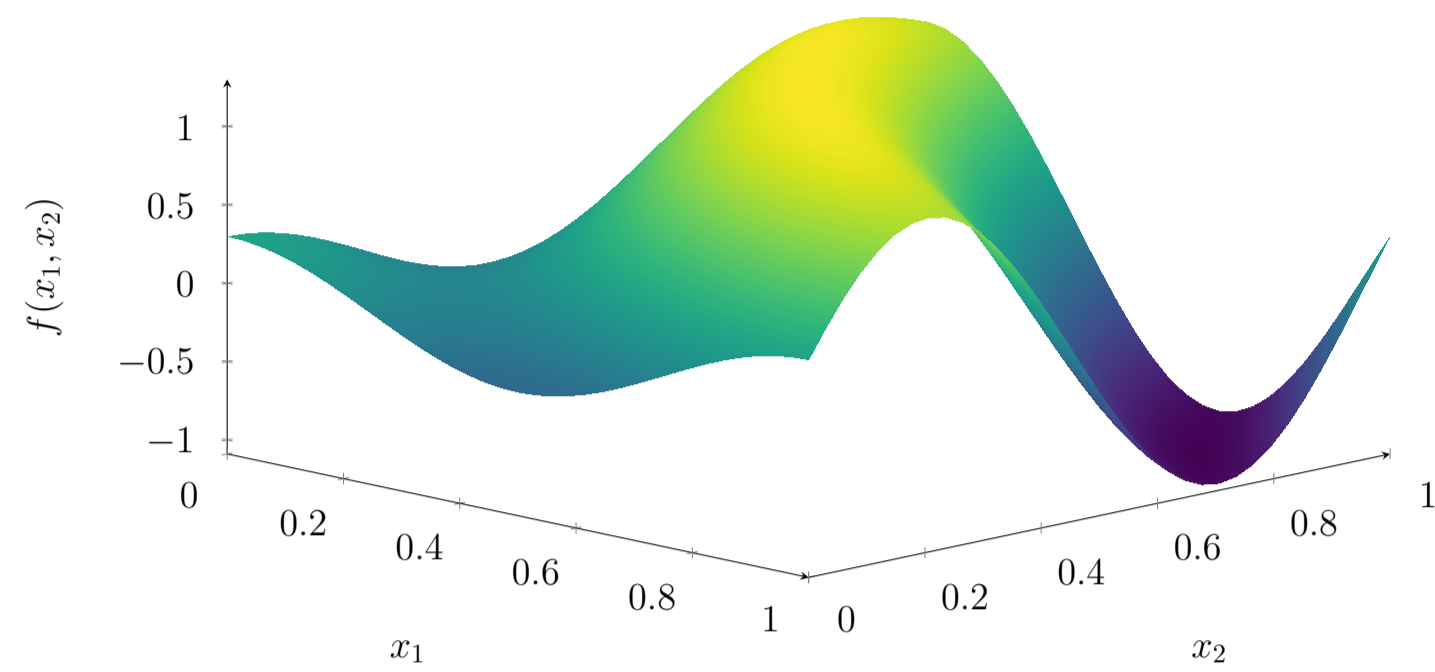
# Optimization, Generalization and Differential Privacy Bounds for Gradient Descent on Kolmogorov–Arnold Networks

Puyu Wang<sup>1</sup>, Junyu Zhou<sup>2</sup>, Philipp Liznerski<sup>1</sup> and Marius Kloft<sup>1</sup>

<sup>1</sup> RPTU Kaiserslautern-Landau, Germany <sup>2</sup> Catholic University of Eichstätt-Ingolstadt, Germany

## Kolmogorov–Arnold Representation: Intuition

**Question.** Can such a multivariate function be represented *exactly* as a **finite sum** of compositions of *univariate* functions and addition?



Example on  $[0, 1]^2$ :  $f(x_1, x_2) = \sin(2\pi x_1 x_2) + 0.3 \cos(2\pi(x_1 + x_2))$ .

**At first glance:** this surface does not look like a “sum of 1D pieces”.

**However,** Kolmogorov–Arnold representation theorem says **yes**, even though the structure is not evident!

**Kolmogorov–Arnold representation theorem (Informal statement).** For any continuous function  $f : [0, 1]^d \rightarrow \mathbb{R}$ , there exist continuous *univariate* functions  $\phi_{q,p} : \mathbb{R} \rightarrow \mathbb{R}$  and  $\Phi_q : \mathbb{R} \rightarrow \mathbb{R}$  such that

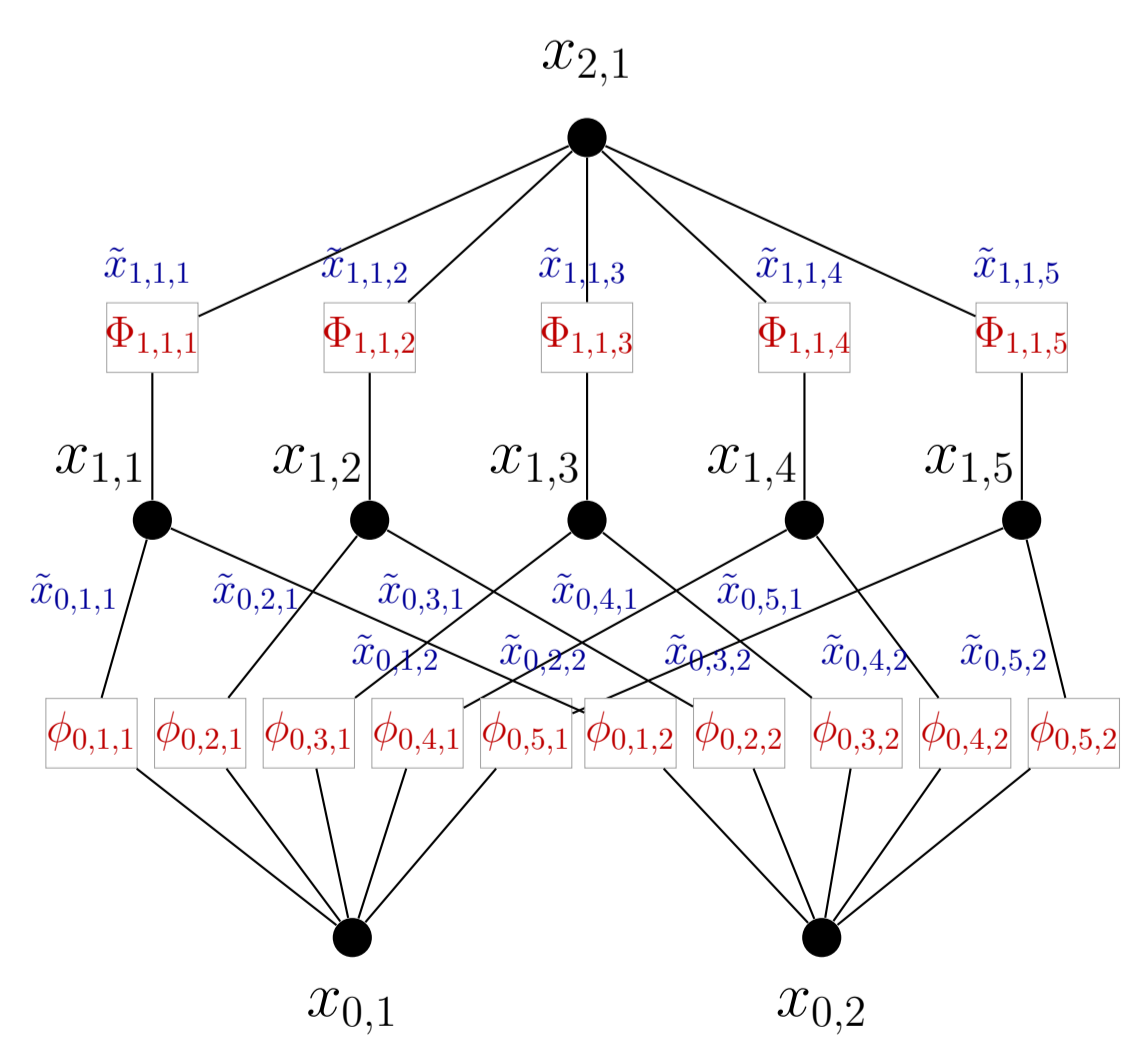
$$f(x_1, \dots, x_d) = \sum_{q=1}^{2d+1} \Phi_q \left( \sum_{p=1}^d \phi_{q,p}(x_p) \right).$$

In other words, any continuous multivariate function can be represented exactly as a **finite sum of compositions of univariate functions and addition**.

**Special case ( $d = 2$ ).** There exist univariate functions  $\phi_{q,1}, \phi_{q,2}$  and  $\Phi_q$  such that

$$f(x_1, x_2) = \sum_{q=1}^5 \Phi_q(\phi_{q,1}(x_1) + \phi_{q,2}(x_2)).$$

## KAN Architecture & Standard Formulation



**KAN idea:** parameterize the Kolmogorov–Arnold representation using **learnable univariate edge functions**.

For  $d = 2$ , a width  $2d+1 = 5$  construction takes the form

$$f(x_{0,1}, x_{0,2}) = \sum_{q=1}^5 \Phi_{1,1,q}(\phi_{0,q,1}(x_{0,1}) + \phi_{0,q,2}(x_{0,2})).$$

Each hidden node corresponds to one index  $q$ , and each edge carries a univariate function.

**Standard KAN layer.** Let  $\mathbf{x}_0 = \mathbf{x} \in \mathbb{R}^{n_0}$ . For each layer  $\ell = 0, \dots, L-1$ , a KAN maps  $\mathbf{x}_\ell \in \mathbb{R}^{n_\ell}$  to  $\mathbf{x}_{\ell+1} \in \mathbb{R}^{n_{\ell+1}}$  by

$$x_{\ell+1,j} = \sum_{i=1}^{n_\ell} \phi_{\ell,j,i}(x_{\ell,i}), \quad j = 1, \dots, n_{\ell+1}.$$

Here each edge carries a learnable univariate function  $\phi_{\ell,j,i} : \mathbb{R} \rightarrow \mathbb{R}$ .

## Why Theory for KANs?

KANs are empirically promising, but theoretical guarantees for **training dynamics**, **generalization**, and **privacy** remain limited.

**We study GD/DP-GD for KANs:**

- **Optimization:** when does GD reduce the training loss?
- **Generalization:** when do GD iterates perform well on test data?
- **Differential privacy:** what utility can private GD guarantee?

**Key scaling question:** **How do  $m, T$ , and  $n$  trade off?**

## Learning Setup for Two-layer KANs

**Two-layer KAN with B-splines.** For  $\mathbf{x}_0 = \mathbf{x} \in \mathbb{R}^d$  and hidden width  $m$ ,

$$f_{\Theta}(\mathbf{x}) = \frac{1}{\sqrt{m}} \sum_{j=1}^m \sum_{k=1}^p c_{j,k} b_k(x_{1,j}),$$

where

$$x_{1,j} = \sigma \left( \frac{1}{\sqrt{d}} \sum_{i=1}^d \sum_{k=1}^p a_{i,j,k} b_k(x_{0,i}) \right), \quad j \in [m].$$

Here  $\Theta = (\mathbf{a}, \mathbf{c})$ , and  $\sigma$  is bounded, Lipschitz, and smooth.

**Risks and training.** For a nonnegative convex, smooth, self-bounded loss  $\ell$ ,

$$\mathcal{L}_S(\Theta) = \frac{1}{n} \sum_{i=1}^n \ell(y_i f_{\Theta}(\mathbf{x}_i)), \quad \mathcal{L}(\Theta) = \mathbb{E}_{(\mathbf{x}, y) \sim P} [\ell(y f_{\Theta}(\mathbf{x}))].$$

We train by gradient descent (GD):

$$\Theta^{(t+1)} = \Theta^{(t)} - \eta \nabla \mathcal{L}_S(\Theta^{(t)}).$$

**We prove upper bounds on optimization risk, population risk, and DP utility,** leading to scaling insights for width  $m$ , iterations  $T$ , and sample size  $n$ .

## Main Results for GD

**NTK separability.** Assume the data are linearly separable with margin  $\gamma > 0$  in the tangent feature space at initialization, i.e., there exists a unit vector  $\Theta_0$ ,

$$y_i \langle \nabla_{\Theta} f_{\Theta_0}(\mathbf{x}_i), \Theta_0 \rangle \geq \gamma, \quad i \in [n].$$

**Optimization.** If  $m \geq \text{polylog}(n)$ , then with high probability over initialization, GD achieves

$$\mathcal{L}_S(\Theta^{(T)}) = \tilde{O}\left(\frac{1}{\gamma^2 \eta T}\right).$$

**Generalization.** If  $m \geq \text{polylog}(n)$  and  $T \gtrsim n$ , then

$$\mathbb{E}_S [\mathcal{L}(\Theta^{(T)})] = \tilde{O}\left(\frac{1}{\gamma^4 n}\right).$$

**Polylogarithmic width already suffices** for  $\tilde{O}(1/T)$  optimization and fast  $\tilde{O}(1/n)$  population risk.

## Differentially Private KANs

**Differential privacy.** For neighboring datasets  $S, S'$  differing in one example, an algorithm  $A$  is  $(\epsilon, \delta)$ -DP if for any event  $E$ ,

$$\mathbb{P}(A(S) \in E) \leq e^\epsilon \mathbb{P}(A(S') \in E) + \delta.$$

**DP-GD.** To achieve  $(\epsilon, \delta)$ -DP, we add Gaussian noise to the gradient and project the iterates back to a neighborhood of initialization:

$$\tilde{\Theta}^{(k+1)} = \Pi_{\Omega} \left( \tilde{\Theta}^{(k)} - \eta (\nabla \mathcal{L}_S(\tilde{\Theta}^{(k)}) + \mathbf{b}^{(k)}) \right), \quad \mathbf{b}^{(k)} \sim \mathcal{N}(0, \sigma^2 I).$$

Here  $\sigma^2$  is calibrated according to  $n, T, \epsilon, \delta$ .

**Private utility.** If  $m \asymp \text{polylog}(n)$  and  $\eta T \asymp \frac{\gamma^2 n \epsilon}{\sqrt{d}}$ , then the  $(\epsilon, \delta)$ -DP variant satisfies

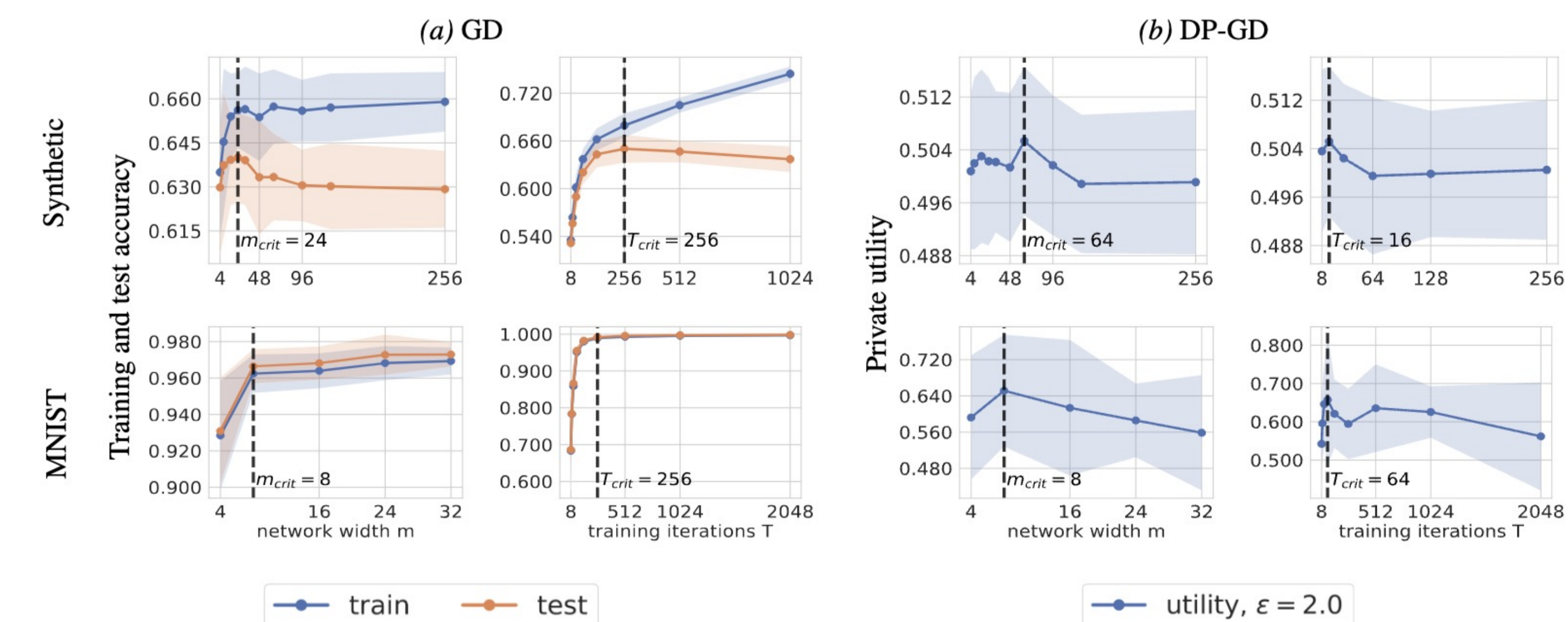
$$\frac{1}{T} \sum_{k=1}^T \mathbb{E}_{S,A} [\mathcal{L}(\tilde{\Theta}^{(k)})] = \tilde{O}\left(\frac{\sqrt{d}}{\gamma^4 n \epsilon}\right).$$

- **Polylogarithmic width is necessary** to achieve  $\tilde{O}\left(\frac{\sqrt{d}}{n \epsilon}\right)$  population risk rate
- **DP-GD attains near-optimal private utility** up to logarithmic factors, despite the nonconvex KAN parameterization

## Practical Implications & Experiments

Our theory gives explicit guidance for choosing the width  $m$  and the number of GD iterations  $T$  for both GD and DP-GD.

- **Training and test accuracy saturate beyond a critical width.**
- **Training accuracy may continue to improve with iterations.**
- **Test accuracy saturates with increasing iterations.**
- **Private utility is maximized in an admissible regime.**



**Figure. Effect of width  $m$  and iterations  $T$  on GD accuracy and DP-GD utility.** Dashed lines indicate empirical change points: accuracy saturates beyond a critical width/training horizon, while DP utility degrades when privacy noise accumulates. These trends match the theory-predicted admissible regimes for choosing  $m$  and  $T$ .