

Super-KAN (Sparse Superposition KAN)

Gary Nan Tie, Oct 26, 2024

Abstract

We introduce a sparse Kolmogorov-Arnold network with superposition layers; each direct summand a version of the Kolmogorov-Arnold representation theorem. Super-KAN are interpretable, enjoy compositional sparsity, and have linear runtime.

Consider a Kolmogorov-Arnold network [Liu et al, 2024]

of shape $[n_0, n_1, \dots, n_L]$

$$\begin{aligned} x_{L+1} &= \Phi_L x_L \\ &= [\varphi_{L,j,i}] [x_{L,i}] \\ &\quad n_{L+1} \times n_L \quad n_L \times 1 \end{aligned}$$

$$\Phi_L : \mathbb{R}^{n_L} \rightarrow \mathbb{R}^{n_{L+1}}, \quad \Phi_L = (f_1, \dots, f_{n_{L+1}}), \quad f_k : \mathbb{R}^{n_L} \rightarrow \mathbb{R}$$

For now, suppress the subscript k as being understood.

$$\text{Let } d = n_L, \quad d \geq 2 \quad \text{and} \quad x = \begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix}.$$

For Super-KAN, define $f : \mathbb{R}^d \rightarrow \mathbb{R}$

as a superposition of shape $[d, 2d+1, 1]$,

$$f(x) \triangleq \Gamma(g) \circ \mathcal{U}(h) x$$

$$\quad 1 \times (2d+1) \quad (2d+1) \times d \quad d \times 1$$

$$\text{where } \mathcal{U}(h) = [\mathcal{U}_{q,p}], \quad \mathcal{U}_{q,p}(x_p) = b_p h(x_p + qa) + c_q$$

For $p = 1, \dots, d$ and $q = 0, \dots, 2d$

and $\Gamma(g) = [g, \dots, g]$
 $1 \times (2d+1)$

For each layer L , and subscript k , we learn:

① univariate functions $g, h: \mathbb{R} \rightarrow \mathbb{R}$,

like wavelets with three parameters (shift, scale, normalization)

② parameters a, b_p, c_q

$$\text{in } f\left(\begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix}\right) = \sum_{q=0}^{2d+1} g\left(\sum_{p=1}^d b_p h(x_p + qa) + c_q\right),$$

a [Braun, 2009] superposition.

Note that f has 2 univariate functions*

and $3d+2$ parameters to learn.

So Super-KAN layer Φ_L has $2n_{L+1}$ univariate

functions and $(3n_L+2)n_{L+1}$ parameters to learn.

In summary, a Super-KAN is a Kolmogorov-Arnold network of shape $[n_0, n_1, \dots, n_L]$, where

$$L\text{-th layer } \Phi_L : \mathbb{R}^{n_L} \rightarrow \mathbb{R}^{n_{L+1}},$$

$$\Phi_L = (f_1^L, \dots, f_{n_{L+1}}^L), \quad f_k^L : \mathbb{R}^{n_L} \rightarrow \mathbb{R}$$

$$f_k^L(x) \triangleq \Gamma(g_k^L) \circ \psi(h_k^L) x, \text{ a Brauer (2009) superposition}$$

$$\text{with } \psi(h_k^L) = [\psi_{q,p}^L], \quad \psi_{q,p}^L(x_p) = b_p h_k^L(x_p + q a) + c_q$$

$$\text{and } a, b_p, c_q \text{ depend on } (L, k), \quad h_k^L : \mathbb{R} \rightarrow \mathbb{R},$$

$$\text{and } \Gamma(g_k^L) = [g_k^L, \dots, g_k^L], \quad g_k^L : \mathbb{R} \rightarrow \mathbb{R}.$$

$1 \times (2n_L + 1)$

Super-KAN are interpretable being a Kolmogorov-Arnold network, compositionally sparse by design, and have linear runtime in terms of univariate functions to compute.

□

Appendix: each superposition is a KAN of shape $[d, 2d+1, 1]$

Kolmogorov-Arnold representation theorem (KART) 1957

For any continuous function $f: [0, 1]^d \rightarrow \mathbb{R}$

there exist univariate continuous functions

$g_q: \mathbb{R} \rightarrow \mathbb{R}$ and $\psi_{p,q}: [0, 1] \rightarrow \mathbb{R}$ such that

$$f(x_1, \dots, x_d) = \sum_{q=0}^{2d} g_q \left(\sum_{p=1}^d \psi_{p,q}(x_p) \right).$$

—/

Theorem (Braun, 2009)

Fix $d \geq 2$. There are real numbers a, b_p, c_q

and a continuous and monotone $\psi: \mathbb{R} \rightarrow \mathbb{R}$

such that for any continuous function $f: [0, 1]^d \rightarrow \mathbb{R}$

there exists a continuous $g: \mathbb{R} \rightarrow \mathbb{R}$ with

$$f(x_1, \dots, x_d) = \sum_{q=0}^{2d} g \left(\sum_{p=1}^d b_p \psi(x_p + qa) + c_q \right).$$

—/

References

Ziming Liu, Yixuan Wang, Sachin Vaidya, Fabian Ruehle,
James Halverson, Marin Soljačić, Thomas Y. Hou,
Max Tegmark. (2024)
KAN: Kolmogorov–Arnold Networks
arXiv:2404.19756v4 [cs.LG]

Kolmogorov, A. N. (1957).
On the representation of continuous functions of many variables
by superposition of continuous functions of one variable and addition.
Doklady Akademii Nauk SSSR, 114, 953–956.

Braun, J. (2009).
An application of Kolmogorov's superposition theorem
to function reconstruction in higher dimensions.
(Ph.D. thesis), Universität Bonn.

Footnote:

* In $f(g, h)$ note that inner h is independent of
target f , ~~outer~~ only outer g depends on f .

So we can further sparsify, to reduce the time to
train and run, by choosing a common inner function
 h across superpositions f in a layer, at the expense
of expressivity.