

Kolmogorov-Arnold Network (KAN)

关键词

取代MLP?

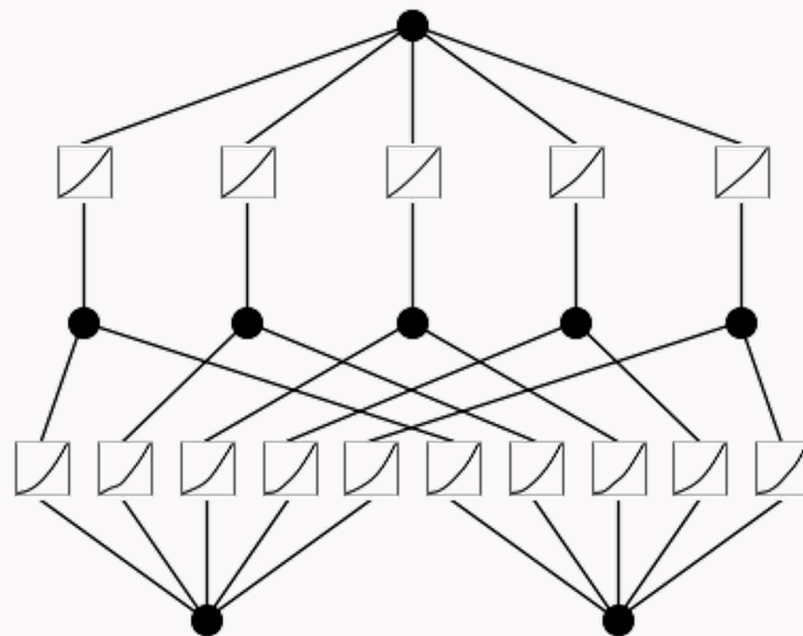
精度更高，参数更少

可解释性

可训练的激活函数

有望解决灾难性遗忘问题

训练时间更长，训练难度更高



Universal Approximation Theorem (通用近似定理)

使 $C(X, \mathbb{R}^m)$ 表示一个从 $X \in \mathbb{R}^n$ 到 \mathbb{R}^m 连续函数集合, 令 $\sigma \in C(\mathbb{R}, \mathbb{R})$, 且 $(\sigma \circ x)_i = \sigma(x_i)$, 即 $\sigma \circ x$ 表示将 σ 应用于 x 的每个分量。

那么当 σ 非多项式时, 当且仅当对于所有 $n \in \mathbb{N}, m \in \mathbb{N}, K$ 是在 \mathbb{R}^n 上的紧子集, $f \in C(K, \mathbb{R}^m), \varepsilon > 0$, 存在 $k \in \mathbb{N}, W \in \mathbb{R}^{k \times n}, b \in \mathbb{R}^k, C \in \mathbb{R}^{k \times m}$, 使得:

$$\sup_{x \in K} \|f(x) - g(x)\| < \varepsilon$$

其中 $g(x) = C \cdot (\sigma \circ (W \cdot x + b))$



Universal Approximation Theorem (通用近似定理)

使 $C(X, \mathbb{R}^m)$ 表示一个从 $X \in \mathbb{R}^n$ 到 \mathbb{R}^m 连续函数集合,令 $\sigma \in C(\mathbb{R}, \mathbb{R})$, 且 $(\sigma \circ x)_i = \sigma(x_i)$,即 $\sigma \circ x$ 表示将 σ 应用于 x 的每个分量。

那么当 σ 非多项式时,当且仅当对于所有 $n \in \mathbb{N}, m \in \mathbb{N}, K$ 是在 \mathbb{R}^n 上的紧子集, $f \in C(K, \mathbb{R}^m), \varepsilon > 0$,存在 $k \in \mathbb{N}, W \in \mathbb{R}^{k \times n}, b \in \mathbb{R}^k, C \in \mathbb{R}^{k \times m}$,使得:

$$\sup_{x \in K} \|f(x) - g(x)\| < \varepsilon$$

其中 $g(x) = C \cdot (\sigma \circ (W \cdot x + b))$

缺点: 可解释性差, 参数量大



Kolmogorov-Arnold表示定理

如果 f 是在有界域上的多元连续函数, 那么 f 可以写成有限数量的单变量连续函数和二元加法运算的组合。进一步来说, 对于一个平滑的函数 $f: [0, 1]^n \rightarrow \mathbb{R}$,

$$f(x) = f(x_1, \dots, x_n) = \sum_{q=1}^{2n+1} \Phi_q \left(\sum_{p=1}^n \phi_{q,p}(x_p) \right)$$

其中 $\phi_{q,p}: [0, 1] \rightarrow \mathbb{R}$, 以及 $\Phi_q: \mathbb{R} \rightarrow \mathbb{R}$ 。

e.g.

$$f(x, y) = xy = \exp(\log(x+1) + \log(y+1)) - \underbrace{(x+0.5)}_{\text{一元函数}} - \underbrace{(y+0.5)}_{\text{一元函数}}$$

自变量为 $\underbrace{\log(x+1)}_{\text{一元函数}} + \underbrace{\log(y+1)}_{\text{一元函数}}$ 的一元函数

Kolmogorov-Arnold表示定理

如果 f 是在有界域上的多元连续函数，那么 f 可以写成有限数量的单变量连续函数和二元加法运算的组合。
进一步来说，对于一个平滑的函数 $f: [0, 1]^n \rightarrow \mathbb{R}$,

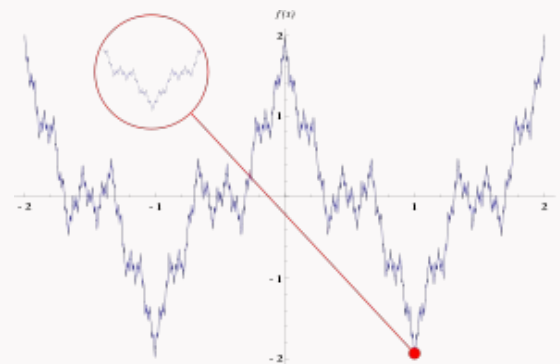
$$f(x) = f(x_1, \dots, x_n) = \sum_{q=1}^{2n+1} \Phi_q \left(\sum_{p=1}^n \phi_{q,p}(x_p) \right)$$

其中内部函数 $\phi_{q,p}: [0, 1] \rightarrow \mathbb{R}$ ，以及外部函数 $\Phi_q: \mathbb{R} \rightarrow \mathbb{R}$ 。

可能是



非平滑函数



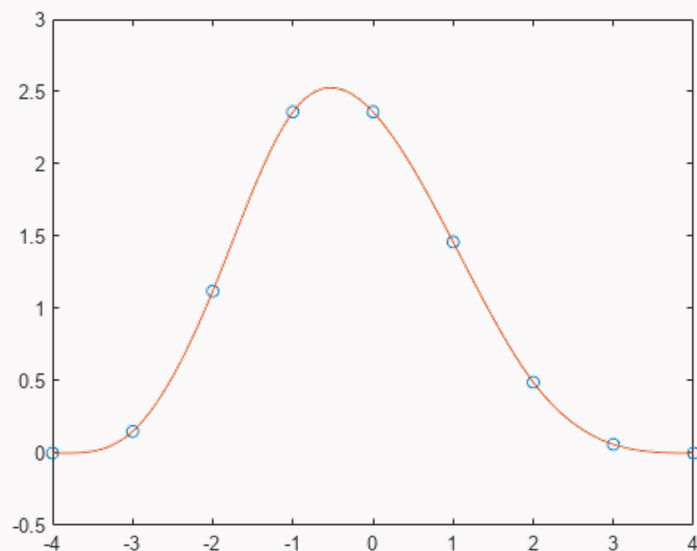
分形函数

不可学习的!

样条函数(Spline)

样条函数是由一些具有连续性条件的子空间上的分段多项式构成, 给定 $n + 1$ 个点 t_0, \dots, t_n 并且满足 $a = t_0 < t_1 < \dots < t_n = b$, 这些点被称为结点(knot), 如果满足下列条件, 参数曲线 $S: [a, b] \rightarrow \mathbb{R}$ 被称为 k 次样条:

1. 在每个分段区间 $[t_i, t_{i+1}]$ 上, S 是一个次数小于等于 k 的多项式。
2. 在 $[t_0, t_n]$ 上 S 有 $k - 1$ 阶连续导数。



优点: 在低维空间中的数据插值和函数逼近中表现出色

缺点: 维数灾难

样条函数(Spline)

样条函数是由一些具有连续性条件的子空间上的分段多项式构成, 给定 $n + 1$ 个点 t_0, \dots, t_n 并且满足 $a = t_0 < t_1 < \dots < t_n = b$, 这些点被称为结点(knot), 如果满足下列条件, 参数曲线 $S: [a, b] \rightarrow \mathbb{R}$ 被称为 k 次样条:

1. 在每个分段区间 $[t_i, t_{i+1}]$ 上, S 是一个次数小于等于 k 的多项式。
2. 在 $[t_0, t_n]$ 上 S 有 $k - 1$ 阶连续导数。

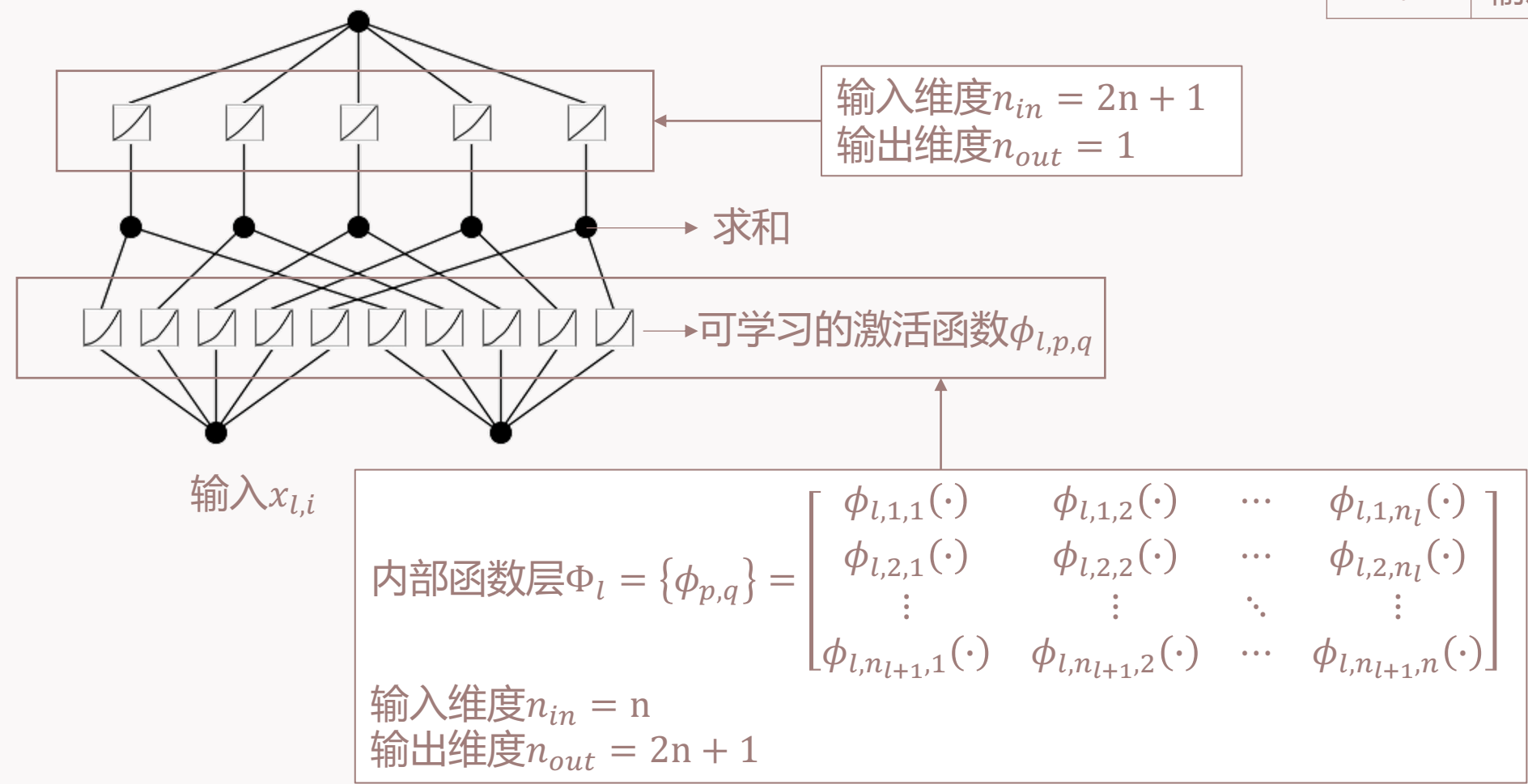
B样条(B-Spline)可以用Cox-de Boor递推公式表达:

$$B_{i,0}(x) := \begin{cases} 1 & \text{if } t_i \leq x < t_{i+1}, \\ 0 & \text{otherwise.} \end{cases}$$
$$B_{i,k}(x) := \frac{x - t_i}{t_{i+k} - t_i} B_{i,k-1}(x) + \frac{t_{i+k+1} - x}{t_{i+k+1} - t_{i+1}} B_{i+1,k-1}(x).$$

KAN架构

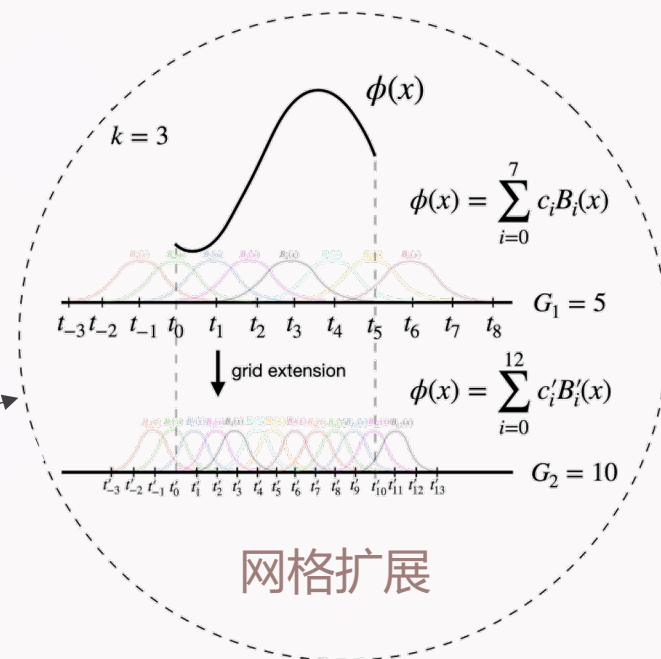
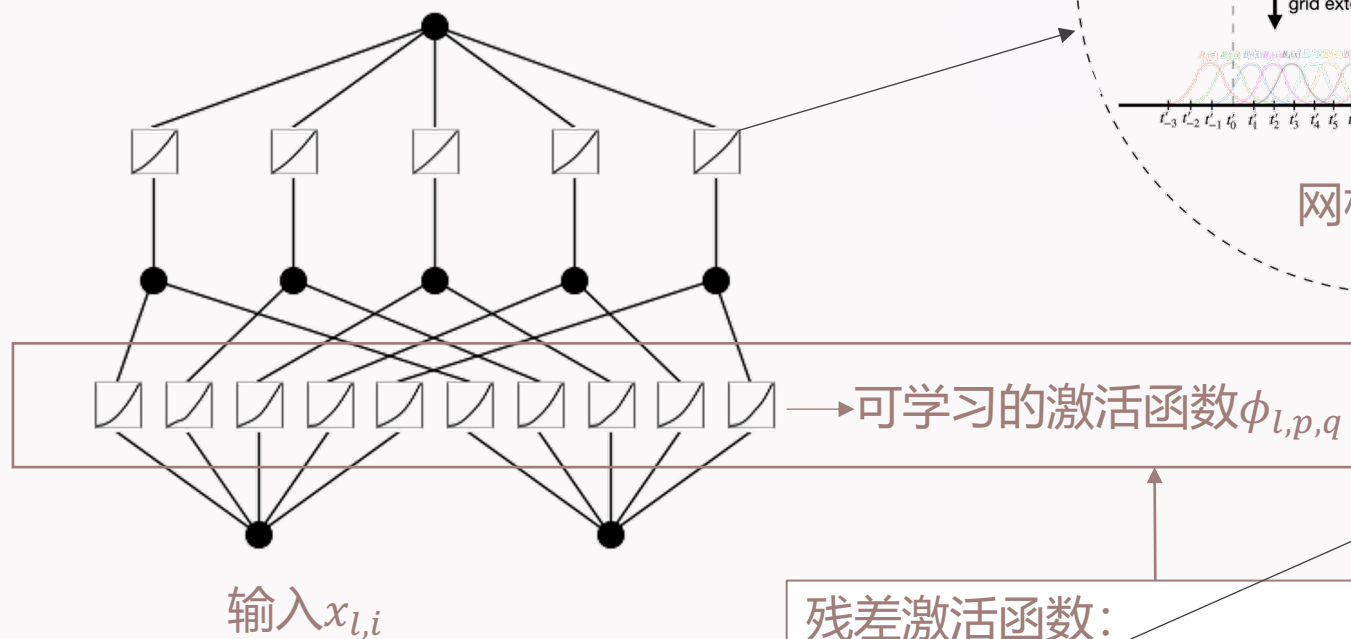
输出: $KAN(x) = (\Phi_{L-1} \circ \Phi_{L-2} \dots \circ \Phi_1 \circ \Phi_0)x$

l	层数的索引
i	x 的索引
p	输入维度 n_{in} 的索引
q	输出维度 n_{out} 的索引
n	输入 x 的维度



KAN架构

输出: $KAN(x) = (\Phi_{L-1} \circ \Phi_{L-2} \dots \circ \Phi_1 \circ \Phi_0)x$



残差激活函数:

$$\phi(x) = w(b(x) + \text{spline}(x))$$

$$b(x) = \text{silu}(x) = \frac{x}{1 + e^{-x}}$$

$\text{spline}(x)$ 被初始化为 $\approx 0^2$

$$\text{spline}(x) = \sum_i c_i B_i(x)$$

w 使用 Xavier 初始化

KAN的可解释性

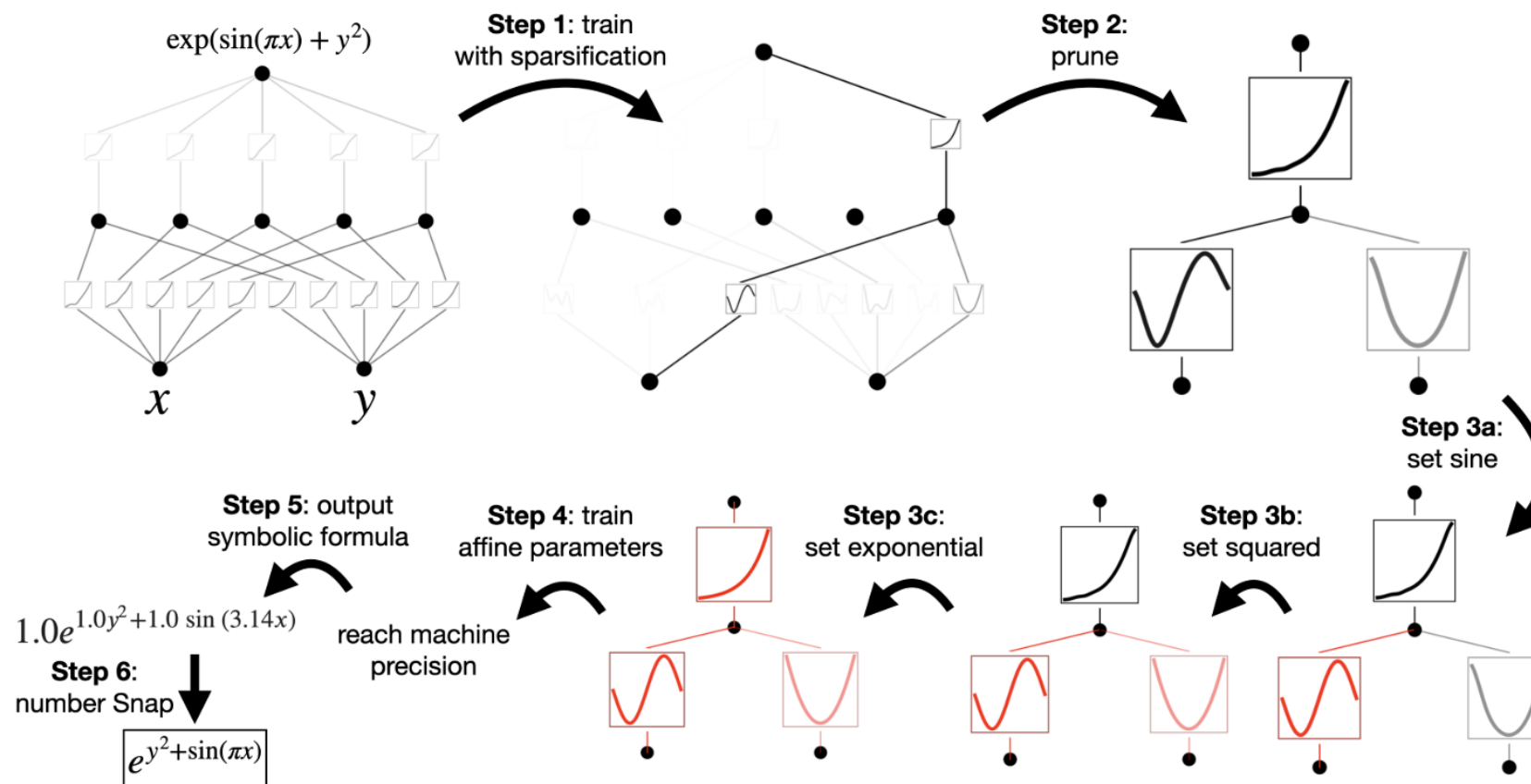


Figure 2.4: An example of how to do symbolic regression with KAN.

拟合特殊函数

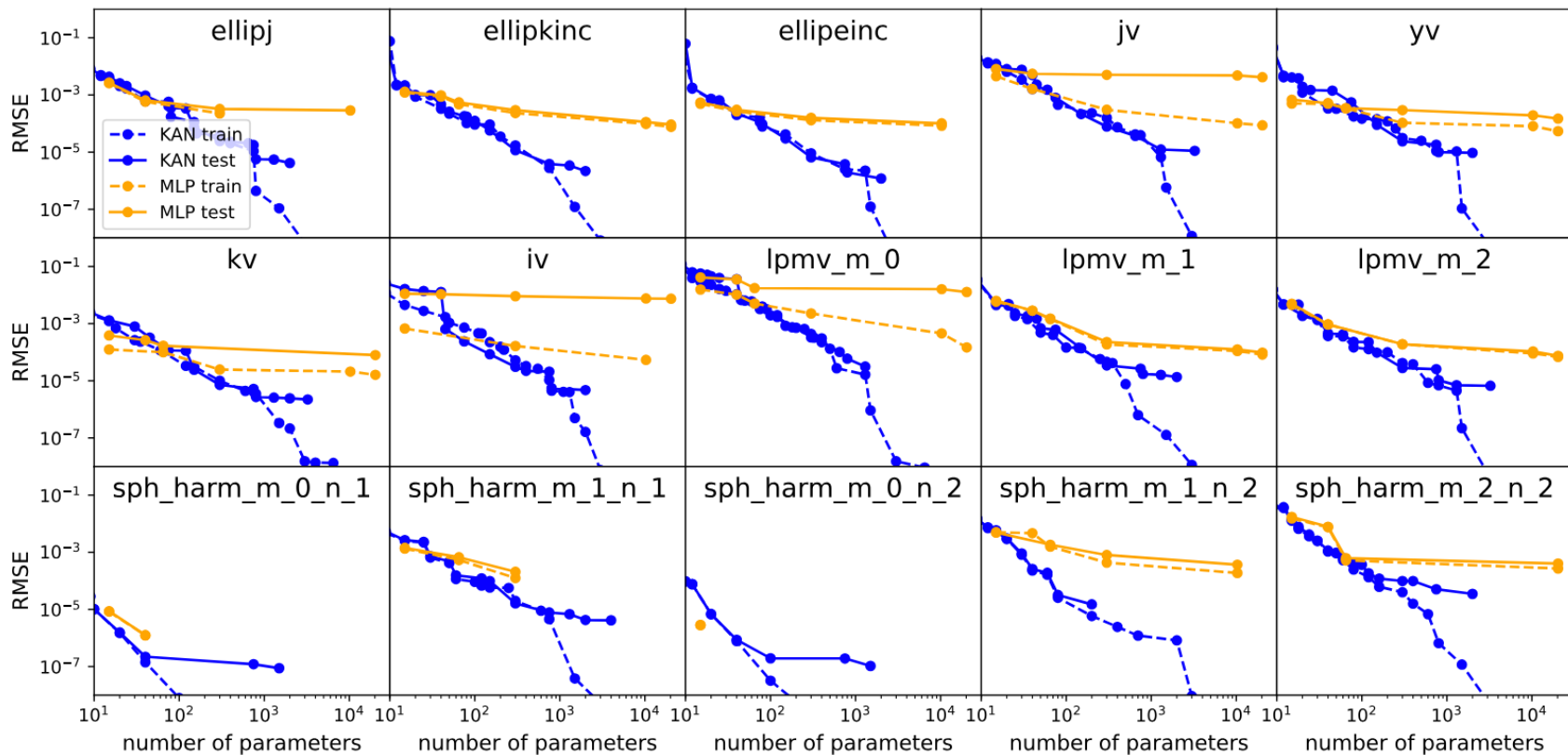


Figure 3.2: Fitting special functions. We show the Pareto Frontier of KANs and MLPs in the plane spanned by the number of model parameters and RMSE loss. Consistently accross all special functions, KANs have better Pareto Frontiers than MLPs. The definitions of these special functions are in Table 2.

解偏微分方程

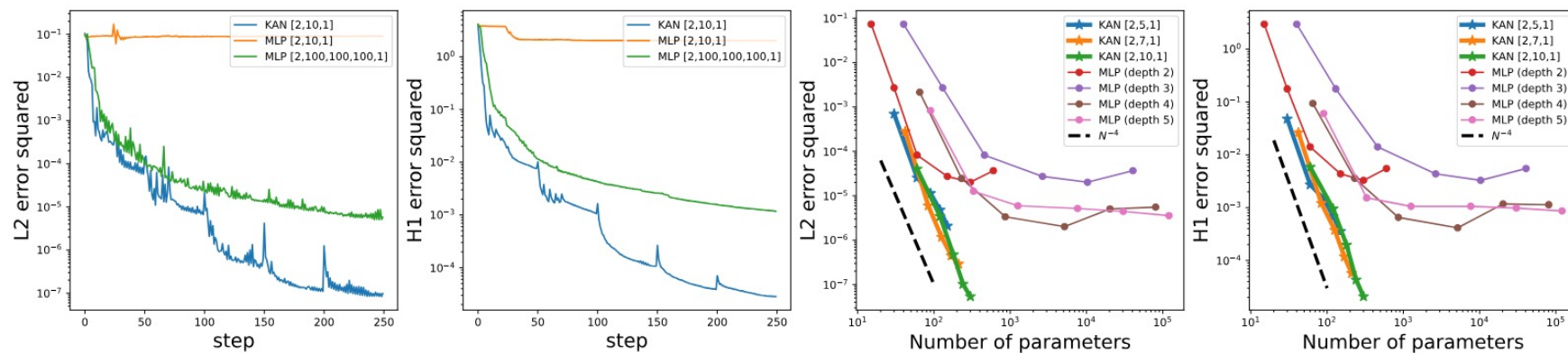


Figure 3.3: The PDE example. We plot L2 squared and H1 squared losses between the predicted solution and ground truth solution. First and second: training dynamics of losses. Third and fourth: scaling laws of losses against the number of parameters. KANs converge faster, achieve lower losses, and have steeper scaling laws than MLPs.

解决灾难性遗忘

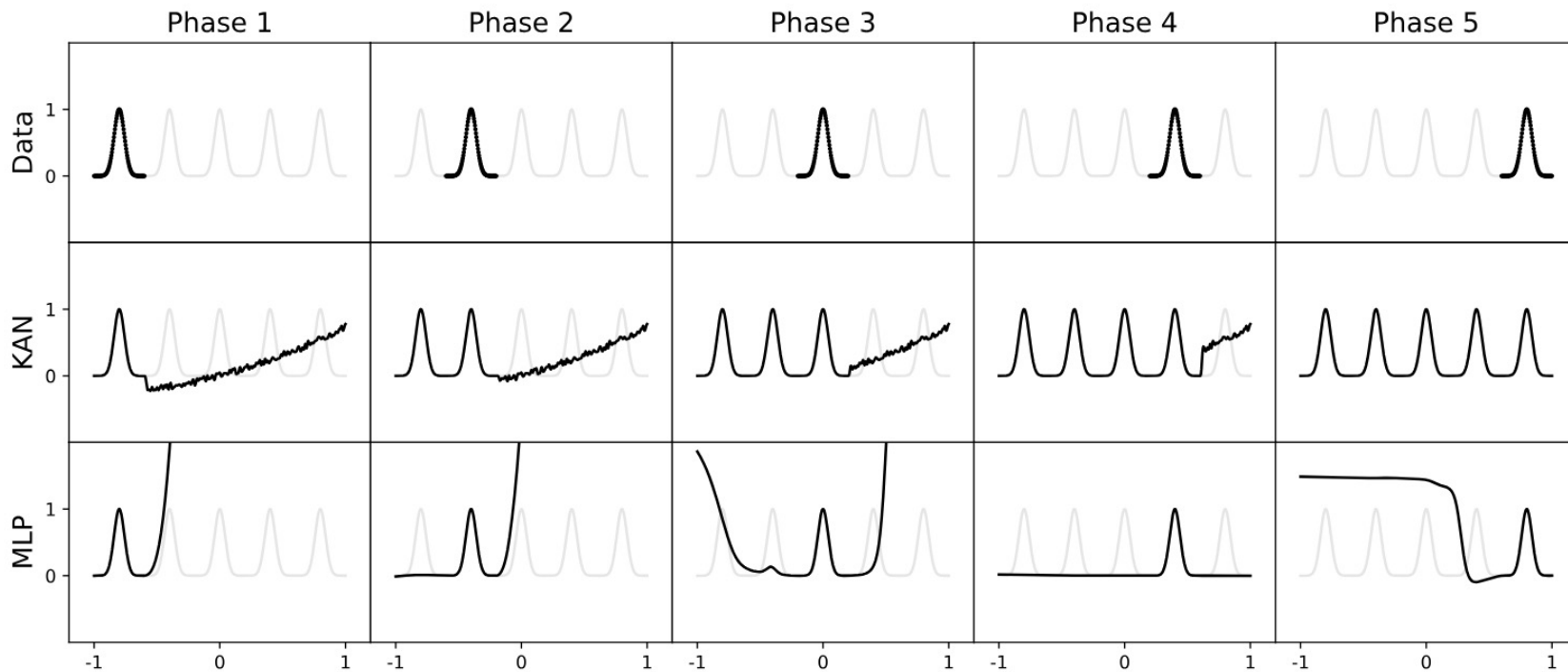


Figure 3.4: A toy continual learning problem. The dataset is a 1D regression task with 5 Gaussian peaks (top row). Data around each peak is presented sequentially (instead of all at once) to KANs and MLPs. KANs (middle row) can perfectly avoid catastrophic forgetting, while MLPs (bottom row) display severe catastrophic forgetting.

MLPs or KANs?

