

# Minimax-optimal nonparametric in-context learning with transformers

Michelle Ching, Ioana Popescu, Nico Smith

Supervisors: Tianyi Ma, William G. Underwood, Richard J. Samworth

Statistical Laboratory, DPMMS

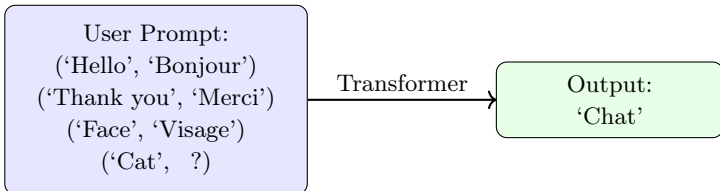
Summer 2025



1. Transformers and their ability to perform regression
2. Why local polynomial regression?
3. Approximation error: Our construction
4. Generalisation error
5. References

Large Language Models (LLMs) such as ChatGPT, LLaMA, and Gemini have shown remarkable in-context learning (ICL) capabilities.

Given a few examples of a task in their prompt, they can often generalise and produce accurate outputs, even if the task was not part of their training set.



This behaviour suggests that LLMs can [learn how to learn](#) from examples in their input.

‘Delivered on time, functions great and long battery life’

Positive

‘My first item was faulty. Was delivered a new one that works great though’

Neutral

‘Not as advertised, customer support was unhelpful when I complained’

Negative

‘Not the best option if you want something to last. Poor battery life also’

?

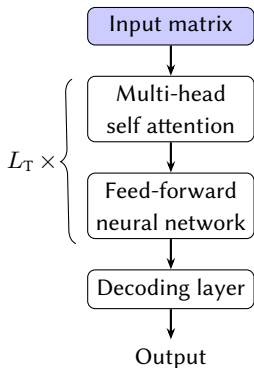
Transformer

Negative

A **transformer** (TF) is an ML architecture that first embeds words as **tokens**, denoted  $\mathbf{h}_1, \dots, \mathbf{h}_{n+1} \in \mathbb{R}^{d_e}$ , and then performs pre-trained operations on them.

These operations consist of **multi-headed self attention** and **feed-forward neural network** (FNN) layers, stacked in  $L_T$  blocks.

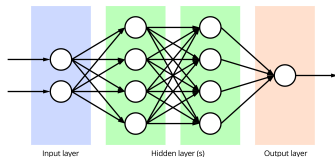
The input to a transformer is  $\mathbf{H} \in \mathbb{R}^{d_e \times (n+1)}$ , with  $i$ th column  $\mathbf{h}_i$ . Combined with the decoding layer, the output is a value in  $\mathbb{R}$ .



- **FNN**: Applies linear transformations and activations **coordinatewise**:

$$\mathbf{h}_i \mapsto \sigma(\mathbf{W}\mathbf{h}_i + \mathbf{b})$$

with weights  $\mathbf{W} \in \mathbb{R}^{d_e \times d_e}$ , bias  $\mathbf{b} \in \mathbb{R}^{d_e}$  and  $\sigma(x) = \text{ReLU}(x) = \max(x, 0)$ .

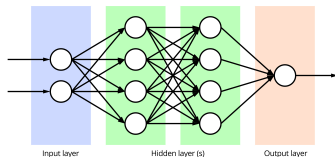


Source\*

- **FNN**: Applies linear transformations and activations **coordinatewise**:

$$\mathbf{h}_i \mapsto \sigma(\mathbf{W}\mathbf{h}_i + \mathbf{b})$$

with weights  $\mathbf{W} \in \mathbb{R}^{d_e \times d_e}$ , bias  $\mathbf{b} \in \mathbb{R}^{d_e}$  and  $\sigma(x) = \text{ReLU}(x) = \max(x, 0)$ .



Source\*

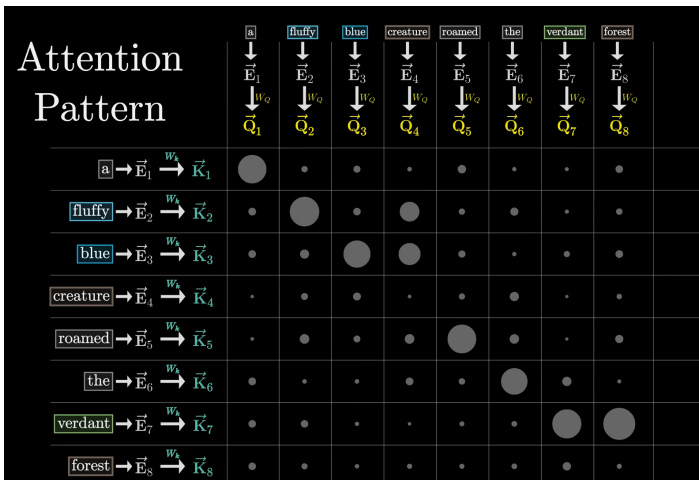
- **Attention**: Can **interact columns**:

$$\mathbf{h}_i \mapsto \sum_{j=1}^n \sigma(\langle \mathbf{Q}\mathbf{h}_i, \mathbf{K}\mathbf{h}_j \rangle) \mathbf{V}\mathbf{h}_j$$

with query  $\mathbf{Q}$ , key  $\mathbf{K}$  and value  $\mathbf{V}$  matrices in  $\mathbb{R}^{d_e \times d_e}$ .

---

\*: <https://dailysciencefactsatit.blogspot.com/2023/06/the-game-changing-inventions-and.html>



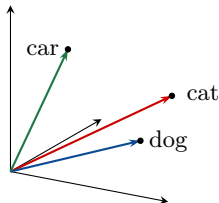
Source: 3Blue1Brown<sup>†</sup>

<sup>†</sup>: Attention in transformers, step-by-step | Deep Learning Chapter 6, YouTube, <https://www.youtube.com/watch?v=eMlx5fFNoYc&t=1065s>



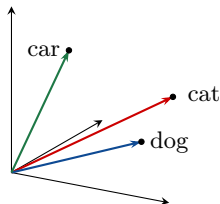
A transformer embeds words as vectors in a  $d_e$ -dimensional space (in practice:  $d_e \approx 10000$ ).

ICL performs **regression** over these embeddings:  
from  $n$  input-output examples, the model infers a  
relation and applies it to a new test point.



A transformer embeds words as vectors in a  $d_e$ -dimensional space (in practice:  $d_e \approx 10000$ ).

ICL performs **regression** over these embeddings:  
from  $n$  input-output examples, the model infers a  
relation and applies it to a new test point.



**Goal:** Analyse the **theoretical** capabilities of transformers to perform **regression**.

**Limitations:** We work with tabular data, where **regression** is more sensible and well-studied – input is in  $\mathbb{R}^d$  with  $d \approx 10$ , and output is in  $\mathbb{R}$ .



**Set-up:** Let  $f_X$  be a density on  $\mathbb{R}^d$  and  $P_{\mathcal{H}}$  be a distribution on the **Hölder class**  $\mathcal{H}(\alpha, L)$  for some  $\alpha, L$ . We consider our data:

- Sample:  $s = (\mathbf{x}_i, y_i)_{i=1}^{n+1}$  i.i.d. pairs where  $\mathbf{x}_i \sim f_X$ ,  $m \sim P_{\mathcal{H}}$ ,  $\varepsilon_i$  is a bounded noise term and  $y_i = m(\mathbf{x}_i) + \varepsilon_i$ .

$$s : \boxed{\boxed{\mathbf{x}_1, y_1}, \boxed{\mathbf{x}_2, y_2}, \dots, \boxed{\mathbf{x}_{n+1}, y_{n+1}}}$$

- Training set:  $\mathcal{S} = (\mathbf{x}_i^{(\gamma)}, y_i^{(\gamma)})_{i \in [n+1], \gamma \in [\Gamma]}$  of independent copies of  $s$ :

$$\mathcal{S} : \boxed{\boxed{s^{(1)}}, \boxed{s^{(2)}}, \dots, \boxed{s^{(\Gamma)}}}$$

Assume  $\|\mathbf{x}_i\|_{\infty} \leq M_x$ ,  $|y_i| \leq M_y$  almost surely.

Input  $\{(\mathbf{x}_i, y_i)_{i=1}^n; \mathbf{x}_{n+1}\}$  is embedded, with **positional encoding**, into:

$$\mathbf{H} = \begin{bmatrix} \mathbf{x}_1 & \cdots & \mathbf{x}_n & \mathbf{x}_{n+1} \\ y_1 & \cdots & y_n & 0 \\ \mathbf{0} & \cdots & \mathbf{0} & \mathbf{0} \\ \mathcal{I}_1 & \cdots & \mathcal{I}_n & \mathcal{I}_{n+1} \\ 1 & \cdots & 1 & 1 \end{bmatrix} \in \mathbb{R}^{d_e \times (n+1)}$$

$$\mathcal{I}_j = \left( \mathbb{1}_{\{j=n+1\}}, \cos\left(\frac{j\pi}{2(n+1)}\right), \sin\left(\frac{j\pi}{2(n+1)}\right) \right)^\top$$

This allows us to distinguish the columns in our constructions.

**Remark.** In practice, LLMs add positional encoding to the word embedding before the first layer of attention.



Denote our **transformer class** by  $\mathcal{T}$ . For  $f \in \mathcal{T}$ , we define

$$\text{Sample loss : } \ell(f, s) = \frac{1}{2} \left( f(\{\mathbf{x}_i, y_i\}_{i=1}^n; \mathbf{x}_{n+1}) - y_{n+1} \right)^2$$

$$\text{Population risk : } R(f) = \mathbb{E}_s[\ell(f, s)]$$

$$\text{Empirical risk : } R_\Gamma(f) = \frac{1}{2\Gamma} \sum_{\gamma=1}^{\Gamma} \left( f(\{\mathbf{x}_i^{(\gamma)}, y_i^{(\gamma)}\}_{i=1}^n; \mathbf{x}_{n+1}^{(\gamma)}) - y_{n+1}^{(\gamma)} \right)^2.$$

Define the empirical minimiser

$$\hat{f} \in \operatorname{argmin}_{f \in \mathcal{T}} R_\Gamma(f).$$



Denote our **transformer class** by  $\mathcal{T}$ . For  $f \in \mathcal{T}$ , we define

$$\text{Sample loss : } \ell(f, s) = \frac{1}{2} \left( f(\{\mathbf{x}_i, y_i\}_{i=1}^n; \mathbf{x}_{n+1}) - y_{n+1} \right)^2$$

$$\text{Population risk : } R(f) = \mathbb{E}_s[\ell(f, s)]$$

$$\text{Empirical risk : } R_\Gamma(f) = \frac{1}{2\Gamma} \sum_{\gamma=1}^{\Gamma} \left( f(\{\mathbf{x}_i^{(\gamma)}, y_i^{(\gamma)}\}_{i=1}^n; \mathbf{x}_{n+1}^{(\gamma)}) - y_{n+1}^{(\gamma)} \right)^2.$$

Define the empirical minimiser

$$\hat{f} \in \operatorname{argmin}_{f \in \mathcal{T}} R_\Gamma(f).$$

**Aim:** Show the empirical risk minimiser  $\hat{f}$ , trained on  $\mathcal{S}$ , can perform well on a new regression task.

**Goal:** Bound  $\mathbb{E}_{\mathcal{S}}[R(\hat{f})] = \mathcal{O}(n^{-\frac{2\alpha}{2\alpha+d}})$ , the optimal non-parametric regression minimax rate.

**Theorem (Kim et al., 2024).** There exists a universal constant  $C > 0$  such that:

$$\mathbb{E}_{\mathcal{S}}[R(\hat{f})] \leq 2 \inf_{f \in \mathcal{T}} R(f) + C \left( \frac{M_y^2}{\Gamma} \log \mathcal{N}(\mathcal{T}, \|\cdot\|_{L^\infty}, \frac{1}{\Gamma}) + \frac{M_y}{\Gamma} \right)$$

where  $\mathcal{N}(\mathcal{T}, \|\cdot\|_{L^\infty}, \delta)$  is the covering number of  $\mathcal{T}$  of radius  $\delta$  with respect to the norm  $\|\cdot\|_{L^\infty}$ .

# Problem set-up

**Goal:** Bound  $\mathbb{E}_{\mathcal{S}}[R(\hat{f})] = \mathcal{O}(n^{-\frac{2\alpha}{2\alpha+d}})$ , the optimal non-parametric regression minimax rate.

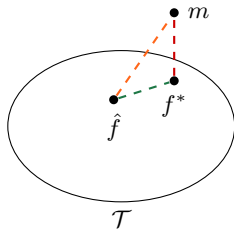
**Theorem (Kim et al., 2024).** There exists a universal constant  $C > 0$  such that:

$$\mathbb{E}_{\mathcal{S}}[R(\hat{f})] \leq 2 \inf_{f \in \mathcal{T}} R(f) + C \left( \frac{M_y^2}{\Gamma} \log \mathcal{N}(\mathcal{T}, \|\cdot\|_{L^\infty}, \frac{1}{\Gamma}) + \frac{M_y}{\Gamma} \right)$$

where  $\mathcal{N}(\mathcal{T}, \|\cdot\|_{L^\infty}, \delta)$  is the covering number of  $\mathcal{T}$  of radius  $\delta$  with respect to the norm  $\|\cdot\|_{L^\infty}$ .

**Approximation error:** How well can a transformer perform? We bound this with an explicit construction  $f^*$ , implementing **local polynomial regression**.

**Generalisation error:** How badly could the empirical risk minimizer  $\hat{f}$  possibly deviate from the true minimizer  $m$ ?





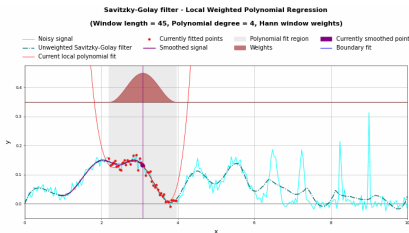
# Local polynomial regression

For  $m \in \mathcal{H}(\alpha, L)$ , local polynomial regression can obtain the minimax optimal pointwise squared error risk of  $\mathcal{O}(n^{-\frac{2\alpha}{2\alpha+d}})$ .

We fit a degree  $p$  polynomial **locally** at  $\mathbf{x}_{n+1}$ , using weighted least squares with **kernel**:

$$K_h(\mathbf{x}_{n+1} - \mathbf{x}_i) := \frac{1}{h^d} K\left(\frac{\mathbf{x}_{n+1} - \mathbf{x}_i}{h}\right)$$

**Note:** We illustrate with  $d = 1$ .



**Figure:** Local weighted polynomial regression (degree=4) smoothing a signal



## Prior works:

- linear and ridge regression (Bai et al., 2023)
- Lasso (Bai et al., 2023)
- Nadaraya–Watson (Shen et al., 2025)
- nonparametric regression with basis expansion (Kim et al., 2024)

Local polynomial regression can be written as a weighted least squares problem, whose solution can be computed by gradient descent.

To fit a polynomial at  $x_{n+1}$ , we must find coefficients for each monomial in  $\mathbf{P}_h(x_{n+1} - x_i)$ , where  $\mathbf{P}(z) = (1, z, z^2, \dots, z^p)^\top$ ,  $\mathbf{P}_h(\cdot) := \mathbf{P}(\frac{\cdot}{h})$ . Let

$$\mathbf{X} := \begin{pmatrix} \mathbf{P}_h(x_{n+1} - x_1)^\top \\ \vdots \\ \mathbf{P}_h(x_{n+1} - x_n)^\top \end{pmatrix} \in \mathbb{R}^{n \times (p+1)} \quad \mathbf{Y} := (y_1 \cdots y_n)^\top \in \mathbb{R}^n$$

$$\mathbf{W} := \text{diag}(K_h(x_{n+1} - x_1), \dots, K_h(x_{n+1} - x_n)) \in \mathbb{R}^{n \times n}$$

To fit a polynomial at  $x_{n+1}$ , we must find coefficients for each monomial in  $\mathbf{P}_h(x_{n+1} - x_i)$ , where  $\mathbf{P}(z) = (1, z, z^2, \dots, z^p)^\top$ ,  $\mathbf{P}_h(\cdot) := \mathbf{P}(\frac{\cdot}{h})$ . Let

$$\mathbf{X} := \begin{pmatrix} \mathbf{P}_h(x_{n+1} - x_1)^\top \\ \vdots \\ \mathbf{P}_h(x_{n+1} - x_n)^\top \end{pmatrix} \in \mathbb{R}^{n \times (p+1)} \quad \mathbf{Y} := (y_1 \cdots y_n)^\top \in \mathbb{R}^n$$

$$\mathbf{W} := \text{diag}(K_h(x_{n+1} - x_1), \dots, K_h(x_{n+1} - x_n)) \in \mathbb{R}^{n \times n}$$

And compute the **weighted least squares estimator**:

$$\begin{aligned} \beta^* &= \underset{\beta \in \mathbb{R}^d}{\text{argmin}} (\mathbf{X}\beta - \mathbf{Y})^\top \mathbf{W} (\mathbf{X}\beta - \mathbf{Y}) \\ &= (\mathbf{X}^\top \mathbf{W} \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{W} \mathbf{Y} \\ &= (\mathbf{X}_W^\top \mathbf{X}_W)^{-1} \mathbf{X}_W^\top \mathbf{Y}_W \end{aligned}$$

where  $\mathbf{X}_W := \sqrt{\mathbf{W}} \mathbf{X}$ ,  $\mathbf{Y}_W := \sqrt{\mathbf{W}} \mathbf{Y}$



We fit  $\beta_0 + \beta_1(x_{n+1} - x_i) + \cdots + \beta_p(x_{n+1} - x_i)^p$  as an expansion around  $x_{n+1}$ , so we approximate the regression function  $m(x_{n+1}) \approx \beta_0$ .

Our construction has 2 steps:



We fit  $\beta_0 + \beta_1(x_{n+1} - x_i) + \cdots + \beta_p(x_{n+1} - x_i)^p$  as an expansion around  $x_{n+1}$ , so we approximate the regression function  $m(x_{n+1}) \approx \beta_0$ .

Our construction has 2 steps:

1: Construct and multiply:

$$\mathbf{P}_h(x_{n+1} - x_i), \sqrt{K_h(x_{n+1} - x_i)} \quad \rightarrow \quad \mathbf{P}_h(x_{n+1} - x_i)\sqrt{K_h(x_{n+1} - x_i)}$$

We fit  $\beta_0 + \beta_1(x_{n+1} - x_i) + \cdots + \beta_p(x_{n+1} - x_i)^p$  as an expansion around  $x_{n+1}$ , so we approximate the regression function  $m(x_{n+1}) \approx \beta_0$ .

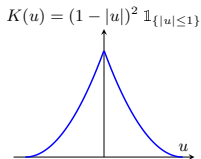
Our construction has 2 steps:

1: Construct and multiply:

$$P_h(x_{n+1} - x_i), \sqrt{K_h(x_{n+1} - x_i)} \quad \rightarrow \quad P_h(x_{n+1} - x_i)\sqrt{K_h(x_{n+1} - x_i)}$$

2: Perform gradient descent using the pairs

$$\left( \sqrt{K_h(x_{n+1} - x_i)} P_h(x_{n+1} - x_i), \sqrt{K_h(x_{n+1} - x_i)} y_i \right)$$



Denoting  $\Delta_j := \mathbf{x}_{n+1} - \mathbf{x}_j$ , column  $j$  undergoes these changes:

$$\begin{bmatrix} \mathbf{x}_j \\ y_j \\ \Delta_j/h \\ \mathbf{0}_{2D+1} \\ \mathcal{I}_j \\ 1 \end{bmatrix} \xrightarrow{\mathbf{F}} \begin{bmatrix} \mathbf{x}_j \\ y_j \\ \Delta_j/h \\ \mathbf{0}_{D-d} \\ \sqrt{K_h(\Delta_j)} \\ \mathbf{0}_{D+1} \\ \mathcal{I}_j \\ 1 \end{bmatrix} \xrightarrow{\mathbf{A}/\mathbf{F}} \begin{bmatrix} \mathbf{x}_j \\ y_j \\ \mathbf{P}_h(\Delta_j) \\ \sqrt{K_h(\Delta_j)} \\ \mathbf{0}_{D+1} \\ \mathcal{I}_j \\ 1 \end{bmatrix} \xrightarrow{\mathbf{A}/\mathbf{F}} \begin{bmatrix} \mathbf{x}_j \\ y_j \\ \mathbf{P}_h(\Delta_j) \\ \sqrt{K_h(\Delta_j)} \\ \sqrt{K_h(\Delta_j)} \mathbf{P}_h(\Delta_j) \\ \sqrt{K_h(\Delta_j)} y_j \\ \mathcal{I}_j \\ 1 \end{bmatrix} \xrightarrow{\mathbf{A}} \begin{bmatrix} \mathbf{x}_j \\ y_j \\ \beta \\ * \end{bmatrix}$$

We then extract  $\beta_0$ . To perform these operations, we need to multiply, copy and add within/between the columns.



## Step 1: Two approaches



$$\begin{bmatrix} \mathbf{x}_j \\ y_j \\ \Delta_j/h \\ \mathbf{0}_{2D+1} \\ \mathcal{I}_j \\ 1 \end{bmatrix} \xrightarrow{\mathbf{F}} \begin{bmatrix} \mathbf{x}_j \\ y_j \\ \Delta_j/h \\ \mathbf{0}_{D-d} \\ \sqrt{K_h(\Delta_j)} \\ \mathbf{0}_{D+1} \\ \mathcal{I}_j \\ 1 \end{bmatrix} \xrightarrow{\mathbf{A}} \begin{bmatrix} \mathbf{x}_j \\ y_j \\ \mathbf{P}_h(\Delta_j) \\ \sqrt{K_h(\Delta_j)} \\ \mathbf{0}_{D+1} \\ \mathcal{I}_j \\ 1 \end{bmatrix} \xrightarrow{\mathbf{A}} \begin{bmatrix} \mathbf{x}_j \\ y_j \\ \mathbf{P}_h(\Delta_j) \\ \sqrt{K_h(\Delta_j)} \\ \sqrt{K_h(\Delta_j)} \mathbf{P}_h(\Delta_j) \\ \sqrt{K_h(\Delta_j)} y_j \\ \mathcal{I}_j \\ 1 \end{bmatrix}$$

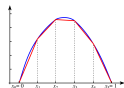
Exact construction (more attention): Allows us to build the monomials exactly.

## Step 1: Two approaches

$$\begin{bmatrix} \mathbf{x}_j \\ y_j \\ \Delta_j/h \\ \mathbf{0}_{2D+1} \\ \mathcal{I}_j \\ 1 \end{bmatrix} \xrightarrow{\text{F}} \begin{bmatrix} \mathbf{x}_j \\ y_j \\ \Delta_j/h \\ \mathbf{0}_{D-d} \\ \sqrt{K_h(\Delta_j)} \\ \mathbf{0}_{D+1} \\ \mathcal{I}_j \\ 1 \end{bmatrix} \xrightarrow{\text{F}} \begin{bmatrix} \mathbf{x}_j \\ y_j \\ P_h(\Delta_j) + \delta_j \\ \sqrt{K_h(\Delta_j)} \\ \mathbf{0}_{D+1} \\ \mathcal{I}_j \\ 1 \end{bmatrix} \xrightarrow{\text{F}} \begin{bmatrix} \mathbf{x}_j \\ y_j \\ P_h(\Delta_j) + \delta_j \\ \sqrt{K_h(\Delta_j)} \\ \sqrt{K_h(\Delta_j)} P_h(\Delta_j) + \epsilon_j \\ \sqrt{K_h(\Delta_j)} y_j + \epsilon_j \\ \mathcal{I}_j \\ 1 \end{bmatrix}$$

Approximate construction (more FNN): Can only approximate the monomials with piecewise linear functions. It is quicker and requires less positional encoding.

**Remark.** A FNN with width  $\mathcal{O}(N)$ , depth  $\mathcal{O}(L)$  can approximate these functions with error  $\mathcal{O}(N^{-L})$ .



**Figure:** Piecewise linear approximation of a quadratic

## Step 2: Performing gradient descent

Attention can perform **1 step of GD in 1 layer of attention** (2 heads).

We rename  $\mathbf{x}_j = \sqrt{K_h(\Delta_j)} \mathbf{P}_h(\Delta_j)$  and  $y_j = \sqrt{K_h(\Delta_j)} y_j$ .

At the  $t^{\text{th}}$  step, each column  $i$  contains  $\mathbf{h}_i = [\mathbf{x}_i, y_i, \beta_{GD}^t, *]^{\top}$  and we want to construct  $\beta_{GD}^{t+1}$ .

$$\dots \xrightarrow{\text{A}} \begin{bmatrix} \mathbf{x}_j \\ y_j \\ \beta_{GD}^t \\ * \end{bmatrix} \xrightarrow{\text{A}} \begin{bmatrix} \mathbf{x}_j \\ y_j \\ \beta_{GD}^{t+1} \\ * \end{bmatrix} \xrightarrow{\text{A}} \begin{bmatrix} \mathbf{x}_j \\ y_j \\ \beta_{GD}^{t+2} \\ * \end{bmatrix} \xrightarrow{\text{A}} \dots$$

We run gradient descent for  $\tau$  steps.

## Step 2: Performing gradient descent

Our loss function

$$\hat{L}_n(\beta) = \frac{1}{2n} \sum_{j=1}^n (\langle \beta, \mathbf{x}_j \rangle - y_j)^2$$

gives update formula:

$$\beta_{GD}^{t+1} = \beta_{GD}^t - \eta \nabla \hat{L}_n(\beta_{GD}^t) = \beta_{GD}^t - \frac{\eta}{n} \sum_{j=1}^n (\langle \beta_{GD}^t, \mathbf{x}_j \rangle - y_j) \mathbf{x}_j$$

We construct the term  $(\langle \beta_{GD}^t, \mathbf{x}_j \rangle - y_j) \mathbf{x}_j$  exactly with attention, using  $u - v = \sigma(u - v) - \sigma(v - u)$ . Recall the form of an attention head:

$$\mathbf{h}_i \mapsto \mathbf{h}_i + A(\mathbf{h}_i) \quad \text{where} \quad A(\mathbf{h}_i) = \sum_{j=1}^n \sigma(\langle \mathbf{Q}\mathbf{h}_i, \mathbf{K}\mathbf{h}_j \rangle) \mathbf{V}\mathbf{h}_j$$

## Step 2: Performing gradient descent

**Remark.** Recall  $\mathbf{X}_W = \sqrt{\mathbf{W}}\mathbf{X}$ . To choose the **step size**  $\eta$ , we look at :

$$\nabla \nabla^\top \hat{L}_n(\boldsymbol{\beta}) = \frac{1}{2n} \nabla \nabla^\top (\mathbf{X}_W \boldsymbol{\beta} - Y_W)^\top (\mathbf{X}_W \boldsymbol{\beta} - Y_W) = \frac{1}{n} \mathbf{X}_W^\top \mathbf{X}_W$$

So choose  $\eta = \frac{1}{\lambda_+} \leq \mathbb{P} \frac{1}{\lambda_{\max}(\mathbf{X}_W^\top \mathbf{X}_W / n)}$ , where  $\lambda_- \leq \lambda_{\min} \leq \lambda_{\max} \leq \lambda_+$ , w.h.p.

## Step 2 (approximate construction)

**Remark.** For the approximate construction using FNNs, we initialise with slightly wrong data  $\tilde{\mathbf{x}}_i, \tilde{y}_i$ , with error bounded by  $\tilde{\varepsilon}$  (in fact data refers to  $\sqrt{K_h(\Delta_i)}P_h(\Delta_i)$  and  $\sqrt{K_h(\Delta_i)}y_i$ ).

This error accumulates at each step since we use the **wrong loss**:

$$\tilde{L}_n(\beta) = \frac{1}{n} \sum_{i=1}^n \ell(\beta^\top \tilde{\mathbf{x}}_i, \tilde{y}_i) = \frac{1}{2n} \sum_{i=1}^n (\beta^\top \tilde{\mathbf{x}}_i - \tilde{y}_i)^2$$

However, after  $t$  steps, it **only accumulates linearly** as  $\tilde{\varepsilon}\eta t$ .



Recall  $m$  is the true regression function,  $\hat{m}_n$  is the local polynomial regression estimator and  $f^*$  is our constructed transformer. We use this to bound the term

$$\inf_{f \in \mathcal{T}} R(f).$$



Recall  $m$  is the true regression function,  $\hat{m}_n$  is the local polynomial regression estimator and  $f^*$  is our constructed transformer. We use this to bound the term

$$\inf_{f \in \mathcal{T}} R(f).$$

**Exact construction:** Choosing bandwidth  $h = h^* = n^{-\frac{1}{2\alpha+d}}$  and defining the noise term  $\varepsilon := y - m$  gives

$$\begin{aligned} R(f^*) &= \mathbb{E} [(f^* - m)^2] + \mathbb{E} [(m - y)^2] \\ &\leq 2\|f^* - \hat{m}_n\|_\infty^2 + 2\mathbb{E} [(\hat{m}_n - m)^2] + \text{Var}(\varepsilon) \\ &= \left[ \mathcal{O}\left(e^{-\frac{\lambda_- \tau}{2\lambda_+}} n^{\frac{d}{2\alpha+d}}\right) + \mathcal{O}\left(n^{-\frac{2\alpha}{2\alpha+d}}\right) + \text{Var}(\varepsilon) \right] \end{aligned}$$

where for  $n$  sufficiently large:  $\lambda_- \leq \lambda_{\min} \leq \lambda_{\max} \leq \lambda_+$ , w.h.p.



## Approximation error (exact construction)

Recall  $m$  is the true regression function,  $\hat{m}_n$  is the local polynomial regression estimator and  $f^*$  is our constructed transformer. We use this to bound the term

$$\inf_{f \in \mathcal{T}} R(f).$$

**Exact construction:** Choosing bandwidth  $h = h^* = n^{-\frac{1}{2\alpha+d}}$  and defining the noise term  $\varepsilon := y - m$  gives

$$\begin{aligned} R(f^*) &= \mathbb{E} [(f^* - m)^2] + \mathbb{E} [(m - y)^2] \\ &\leq 2\|f^* - \hat{m}_n\|_\infty^2 + 2\mathbb{E} [(\hat{m}_n - m)^2] + \text{Var}(\varepsilon) \\ &= \boxed{\mathcal{O}\left(e^{-\frac{\lambda_- \tau}{2\lambda_+}} n^{\frac{d}{2\alpha+d}}\right) + \mathcal{O}\left(n^{-\frac{2\alpha}{2\alpha+d}}\right) + \text{Var}(\varepsilon)} \end{aligned}$$

where for  $n$  sufficiently large:  $\lambda_- \leq \lambda_{\min} \leq \lambda_{\max} \leq \lambda_+$ , w.h.p.

So we require  $\boxed{\tau \geq \frac{2\lambda_+}{\lambda_-} \log n}$  to match minimax rate.



**FNN construction:** Choosing FNNs with width  $\mathcal{O}(N)$ , depth  $\mathcal{O}(L)$ :

$$R(f^*) \leq \mathcal{O}\left(n^{\frac{\alpha+p+d}{2\alpha+d}} N^{-L} \tau\right) \\ + \mathcal{O}\left(e^{-\frac{\lambda_- \tau}{2\lambda_+}} n^{\frac{d}{2\alpha+d}}\right) + \mathcal{O}\left(n^{-\frac{2\alpha}{2\alpha+d}}\right) + \text{Var}(\varepsilon)$$

where the first term comes from **error in data**, accumulated over  $\tau$  steps of gradient descent. So to match the **minimax rate**, we take

$$\tau \geq \frac{2\lambda_+}{\lambda_-} \log n, \quad N = \text{const}, \quad L = \mathcal{O}(\log n).$$



We now want to bound the generalisation error term

$$\frac{M_y^2}{\Gamma} \log \mathcal{N}\left(\mathcal{T}, \|\cdot\|_{L^\infty}, \frac{1}{\Gamma}\right) + \frac{M_y}{\Gamma}$$

where  $|y_i| \leq M_y$  a.s. and  $\Gamma$  is our number of training samples.



We now want to bound the generalisation error term

$$\frac{M_y^2}{\Gamma} \log \mathcal{N}\left(\mathcal{T}, \|\cdot\|_{L^\infty}, \frac{1}{\Gamma}\right) + \frac{M_y}{\Gamma}$$

where  $|y_i| \leq M_y$  a.s. and  $\Gamma$  is our number of training samples.

The **covering number**  $\mathcal{N}$  for our transformer class  $\mathcal{T}$  is well studied (Bai et al., 2023). It depends on width, depth, number of attention heads, parameter size, and embedding dimension.

These parameter bounds give us a **Lipschitz constant** for  $\mathcal{T}$  so we can construct an  $\|\cdot\|_{L^\infty}$  net for  $\mathcal{T}$  of a given size.



$$\frac{M_y^2}{\Gamma} \log \mathcal{N}\left(\mathcal{T}, \|\cdot\|_{L^\infty}, \frac{1}{\Gamma}\right) + \frac{M_y}{\Gamma}$$

For  $\tau$  steps of gradient descent:

Exact construction:

$$\log \mathcal{N}\left(\mathcal{T}, \|\cdot\|_{L^\infty}, \frac{1}{\Gamma}\right) \leq \mathcal{O}\left(\tau^2 \log(\Gamma \tau)\right)$$

$$\text{Exact Generalisation Error} \leq \mathcal{O}\left(\frac{\tau^2 \log(\Gamma \tau)}{\Gamma}\right)$$

$$\frac{M_y^2}{\Gamma} \log \mathcal{N}\left(\mathcal{T}, \|\cdot\|_{L^\infty}, \frac{1}{\Gamma}\right) + \frac{M_y}{\Gamma}$$

For  $\tau$  steps of gradient descent:

Exact construction:

$$\log \mathcal{N}\left(\mathcal{T}, \|\cdot\|_{L^\infty}, \frac{1}{\Gamma}\right) \leq \mathcal{O}\left(\tau^2 \log(\Gamma \tau)\right)$$

$$\text{Exact Generalisation Error} \leq \mathcal{O}\left(\frac{\tau^2 \log(\Gamma \tau)}{\Gamma}\right)$$

FNN construction: width  $\mathcal{O}(N)$ , depth  $\mathcal{O}(L)$

$$\log \mathcal{N}\left(\mathcal{T}, \|\cdot\|_{L^\infty}, \frac{1}{\Gamma}\right) \leq \mathcal{O}\left(N^2(\tau + L)^2 \log(\Gamma N(\tau + L))\right)$$

$$\text{FNN Generalisation Error} \leq \mathcal{O}\left(\frac{N^2(\tau + L)^2 \log(\Gamma N(\tau + L))}{\Gamma}\right)$$

$$\text{Exact Generalisation Error} \leq \mathcal{O}\left(\frac{\tau^2 \log(\Gamma \tau)}{\Gamma}\right)$$

$$\text{FNN Generalisation Error} \leq \mathcal{O}\left(\frac{N^2(\tau + L)^2 \log(\Gamma N(\tau + L))}{\Gamma}\right)$$

Recall we earlier took

$$\tau \geq \frac{2\lambda_+}{\lambda_-} \log(n), \quad N = \text{const}, \quad L = \mathcal{O}(\log(n))$$

so to match the minimax rate of  $\mathcal{O}(n^{-\frac{2\alpha}{2\alpha+d}})$ , we need to take:

$$\text{Exact:} \quad \Gamma \geq \mathcal{O}\left(n^{\frac{2\alpha}{2\alpha+d}} (\log n)^3\right)$$

$$\text{FNN:} \quad \Gamma \geq \mathcal{O}\left(n^{\frac{2\alpha}{2\alpha+d}} (\log n)^3\right)$$



**Lemma.** In the exact set up with  $h = h^* = n^{-\frac{1}{2\alpha+d}}$ ,  $\tau \geq \frac{2\lambda_+}{\lambda_-} \log(n)$ , and  $\Gamma \geq \mathcal{O}(n^{\frac{2\alpha}{2\alpha+d}} \log(n)^3)$ , we achieve the non-parametric minimax rate

$$\mathbb{E}_{\mathcal{S}}[R(\hat{f})] \leq \mathcal{O}(n^{-\frac{2\alpha}{2\alpha+d}}) + \text{Var}(\varepsilon)$$





**Lemma.** In the exact set up with  $h = h^* = n^{-\frac{1}{2\alpha+d}}$ ,  $\tau \geq \frac{2\lambda_+}{\lambda_-} \log(n)$ , and  $\Gamma \geq \mathcal{O}(n^{\frac{2\alpha}{2\alpha+d}} \log(n)^3)$ , we achieve the non-parametric minimax rate

$$\mathbb{E}_{\mathcal{S}}[R(\hat{f})] \leq \mathcal{O}(n^{-\frac{2\alpha}{2\alpha+d}}) + \text{Var}(\varepsilon)$$

**Lemma.** In the FNN set up with  $h = h^* = n^{-\frac{1}{2\alpha+d}}$ ,  $\tau \geq \frac{2\lambda_+}{\lambda_-} \log(n)$ ,  $\Gamma \geq \mathcal{O}(n^{\frac{2\alpha}{2\alpha+d}} \log(n)^3)$ , and constant  $N$ ,  $L = \mathcal{O}(\log n)$ , we achieve the non-parametric minimax rate

$$\mathbb{E}_{\mathcal{S}}[R(\hat{f})] \leq \mathcal{O}(n^{-\frac{2\alpha}{2\alpha+d}}) + \text{Var}(\varepsilon)$$



## Achievements:

1. We obtain the optimal non-parametric regression minimax rate of  $\mathcal{O}(n^{-\frac{2\alpha}{2\alpha+d}})$  of pointwise squared error risk in both the FNN construction and exact construction.



## Achievements:

1. We obtain the optimal non-parametric regression minimax rate of  $\mathcal{O}(n^{-\frac{2\alpha}{2\alpha+d}})$  of pointwise squared error risk in both the FNN construction and exact construction.
2.  $\Gamma$ , number of training samples, is controlled by  $\mathcal{O}(n^{\frac{2\alpha}{2\alpha+d}} \log(n)^3)$ , which is less than previous requirements in the literature.



- Y. Bai, F. Chen, H. Wang, C. Xiong, and S. Mei. Transformers as statisticians: Provable in-context learning with in-context algorithm selection, 2023. URL <https://arxiv.org/abs/2306.04637>.
- J. Kim, T. Nakamaki, and T. Suzuki. Transformers are minimax optimal nonparametric in-context learners. In The Thirty-eighth Annual Conference on Neural Information Processing Systems, 2024. URL <https://openreview.net/forum?id=hF6vatntqc>.
- Z. Shen, A. Hsu, R. Lai, and W. Liao. Understanding in-context learning on structured manifolds: Bridging attention to kernel methods, 2025. URL <https://arxiv.org/abs/2506.10959>.