

TRANSFORMERS HANDLE ENDOGENEITY IN IN-CONTEXT LINEAR REGRESSION

Haodong Liang

UC Davis

hdliang@ucdavis.edu

Krishnakumar Balasubramanian

UC Davis

kbala@ucdavis.edu

Lifeng Lai

UC Davis

lflai@ucdavis.edu

ABSTRACT

We explore the capability of transformers to address endogeneity in in-context linear regression. Our main finding is that transformers inherently possess a mechanism to handle endogeneity effectively using instrumental variables (IV). First, we demonstrate that the transformer architecture can emulate a gradient-based bi-level optimization procedure that converges to the widely used two-stage least squares (2SLS) solution at an exponential rate. Next, we propose an in-context pretraining scheme and provide theoretical guarantees showing that the global minimizer of the pre-training loss achieves a small excess loss. Our extensive experiments validate these theoretical findings, showing that the trained transformer provides more robust and reliable in-context predictions and coefficient estimates than the 2SLS method, in the presence of endogeneity.

1 INTRODUCTION

The transformer architecture (Vaswani et al., 2017) has demonstrated remarkable in-context learning (ICL) capabilities across various domains, such as natural language processing (Devlin et al., 2019; Radford et al., 2019; Brown et al., 2020), computer vision (Dosovitskiy et al., 2021; Carion et al., 2020), and reinforcement learning (Lee et al., 2022; Parisotto et al., 2020). Self-attention mechanism, a core component of transformers, allows these models to capture long-range dependencies in data, which is critical for success in these tasks. Despite their impressive performance, the theoretical understanding of transformers remains limited, leaving important questions unanswered about their true capabilities and the underlying mechanisms driving their exceptional results.

Recent efforts to theoretically understand transformers’ ICL capabilities have focused on their performance in fundamental statistical tasks. Focusing on simple function classes, Garg et al. (2022) highlighted that transformers, when trained on sufficiently large and diverse data from a specific function class, can generalize across most functions of that class without task-specific fine-tuning. Building on this, subsequent work by Bai et al. (2023) established that attention layers enable transformers to perform gradient descent, implementing algorithms like linear regression, logistic regression, and LASSO; see also Akyürek et al. (2023); Von Oswald et al. (2023); Li et al. (2023); Fu et al. (2023); Ahn et al. (2023); Jin et al. (2025). The learning dynamics of transformer was analyzed in Huang et al. (2024). Furthermore Zhang et al. (2024a;b) showed that *trained* transformers’ ICL abilities for linear regression tasks are theoretically robust under certain distributional shifts.

Existing works on analyzing the ICL ability of transformers for linear regression tasks, however, ignore *endogeneity* and have mainly focused on the *exogenous* setup where the additive noise is uncorrelated with the explanatory variables. Ignoring *endogeneity* in linear regression leads to biased and inconsistent estimates, resulting from issues like omitted variable bias, simultaneity, and measurement error, which can distort causal inferences and lead to incorrect policy conclusions (Hausman, 2001; Wooldridge, 2015; Angrist & Pischke, 2009; Greene, 2018). Instrumental variable (IV) regression is a widely adopted method to handle endogeneity by utilizing instruments that are correlated with the endogenous variables but uncorrelated with the error term (Angrist & Krueger, 2001). A naturally intriguing question that therefore arises is:

Can transformers leverage instrumental variables and provide reliable predictions and coefficient estimates, in the presence of endogeneity?

In this work, we aim to answer this question and offer new insights on in-context linear regression tasks. Our key contributions include:

- We demonstrate that looped transformers can address endogeneity in linear regression by leveraging instrumental variables. Specifically, we show that transformers can implement two-stage least squares (2SLS) regression through a bi-level gradient descent procedure, where each iteration is executed by a two-layer transformer block. Moreover, the convergence rate to the 2SLS estimator is exponential with respect to the number of blocks.
- We propose an ICL training scheme for transformers to efficiently handle endogeneity. Under this scheme, we show that the global minimizer of the in-context pre-training loss achieves a small excess loss compared to the global optimal expected loss.
- We evaluate the performance of the trained transformer model through extensive experiments, finding that it not only matches the performance of the 2SLS estimator on standard IV tasks but also generalizes effectively to more complex scenarios, including the challenging cases of weak instruments, non-linear IV, and underdetermined IV problems.
- As part of our analysis, we derive the first non-asymptotic bound for the 2SLS estimator under random design, providing valuable insights for future theoretical work.

1.1 RELATED WORKS

In-context Learning. Initial works by Garg et al. (2022) and Bai et al. (2023) adopted the standard multi-layer transformer architecture to conduct the experiments. Later, Giannou et al. (2023) and Yang et al. (2024) showed that a looped architecture reduces the required depth of transformers and exhibits better efficiency in learning algorithms. Gao et al. (2024) illustrated that the looped transformer architecture with extra pre-processing and post-processing layers can achieve higher expressive power than a standard transformer with the same number of parameters. Apart from works concerning the implementability of first-order gradient descent algorithms by transformers, other works have also examined higher-order and non-parametric optimization methods. Specifically, Giannou et al. (2024) showed that transformers can emulate Newton’s method for logistic regression. Cheng et al. (2024) showed that transformers can implement functional gradient descent and hence enable them to learn non-linear functions in-context. Relationship between in-context learning and Bayesian inference is also studied in Ye et al. (2024); Falck et al. (2024).

Nichani et al. (2024) illustrated how the transformers can learn the causal structure by encoding the latent causal graph in the first attention layer. Goel & Bartlett (2024) explored the representational power of transformer for learning linear dynamical systems. Makkuva et al. (2024a;b); Rajaraman et al. (2024); Edelman et al. (2024) considered ICL Markov chains with transformers, including both landscape and training dynamics analyses. To the best of our knowledge, we are not aware of prior works on handling endogeneity with transformers.

Instrumental Variable Regression. IV regression has been widely studied in econometrics (Angrist & Krueger, 2001; Angrist & Pischke, 2009). Recent works in machine learning explored the optimization based approaches for the IV regression problem. Singh et al. (2019) proposed the kernel IV regression to model non-linear relationship between variables. Muandet et al. (2020) proposed that a non-linear IV regression problem can be formulated as a convex-concave saddle point problem. Della Vecchia & Basu (2023); Chen et al. (2024); Fonseca et al. (2024) proposed a stochastic optimization algorithm for IV regression.

Notation: Throughout this paper, unless otherwise specified, lower-case letters denote random variables or samples, while upper-case letters represent datasets (collections of samples). Bolded letters indicate vectors or matrices, whereas unbolded letters indicate scalars. The notation $\mathbf{X}_{:,i}$ refers to the i -th column, and $\mathbf{X}_{i,:}$ refers to the i -th row of matrix \mathbf{X} . $\lambda_{\min}(\cdot)$ denotes the minimum eigenvalue, and $\sigma_{\min}(\cdot)$ denotes the minimum singular value of a matrix. By default, $\|\cdot\|$ denotes the Euclidean norm for a vector, or the spectral norm for a matrix.

2 ENDOGENEITY AND INSTRUMENTAL VARIABLE REGRESSION

Suppose we are interested in estimating the relationship between response variable $y \in \mathbb{R}$ and predictor variable $\mathbf{x} \in \mathbb{R}^p$ with endogeneity. Given instruments $\mathbf{z} \in \mathbb{R}^q$, we consider the model

$$y = \boldsymbol{\beta}^\top \mathbf{x} + \epsilon_1, \quad \text{and} \quad \mathbf{x} = \boldsymbol{\Theta}^\top \mathbf{z} + \epsilon_2, \quad (1)$$

where $\boldsymbol{\beta} \in \mathbb{R}^p$, and $\boldsymbol{\Theta} \in \mathbb{R}^{q \times p}$ are the true model parameters, $\epsilon_1 \in \mathbb{R}$ and $\epsilon_2 \in \mathbb{R}^p$ are (centered) random noise terms with variance σ_1^2 and covariance matrix $\boldsymbol{\Sigma}_2$, respectively. Further, ϵ_2 is an unobserved noise correlated with ϵ_1 , leading to the correlation between \mathbf{x} and ϵ_1 , which introduces confounding in the model between \mathbf{x} and y . Under this setting, the standard ordinary least squares (OLS) estimator is a biased and inconsistent estimator of $\boldsymbol{\beta}$ (see Wooldridge (2015), Chapter 9). To address this issue, instrumental variable (IV) regression is a widely used method to provide a consistent estimate for $\boldsymbol{\beta}$.

Definition 2.1 (2SLS estimator). IV regression is a regression model to provide consistent estimate on the causal effect $\boldsymbol{\beta}$ for the endogeneity problem (1), by utilizing the instrument \mathbf{z} . Given observational values $(\mathbf{Z}, \mathbf{X}, \mathbf{Y}) = \{(\mathbf{z}_i, \mathbf{x}_i, y_i)\}_{i=1}^n$, the standard approach to estimate the IV regression model is 2SLS; see, for example, Wooldridge (2015), Chapter 15.

i. *First stage*: Regress \mathbf{X} on \mathbf{Z} to obtain $\hat{\boldsymbol{\Theta}}$

$$\hat{\boldsymbol{\Theta}} = (\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top \mathbf{X}.$$

ii. *Second stage*: Regress \mathbf{Y} on $\mathbf{Z}\hat{\boldsymbol{\Theta}}$ to obtain:

$$\hat{\boldsymbol{\beta}}_{2\text{SLS}} = (\hat{\boldsymbol{\Theta}}^\top \mathbf{Z}^\top \mathbf{Z} \hat{\boldsymbol{\Theta}})^{-1} \hat{\boldsymbol{\Theta}}^\top \mathbf{Z}^\top \mathbf{Y}. \quad (2)$$

We introduce the standard assumptions required to show the convergence rate of the above estimator.

Assumption 1 (Instrumental variable). A random variable $\mathbf{z} \in \mathbb{R}^q$ is a valid IV, if it satisfies the following conditions:

- i. Fully identification: $q \geq p$ (without loss of generality, we assume data \mathbf{Z}, \mathbf{X} are full rank).
- ii. Correlated to \mathbf{x} : $\text{Corr}(\mathbf{z}, \mathbf{x}) \neq \mathbf{0}$.
- iii. Conditional uncorrelated to y : $\text{Corr}(\mathbf{z}, \epsilon_1) = 0$.

In particular, condition (i) above ensures the existence of unique solution for $\hat{\boldsymbol{\beta}}_{2\text{SLS}}$. We refer to Stock & Watson (2011, Chapter 12) for additional elaborate discussions on the above conditions. To derive non-asymptotic convergence rates, we further assume the following regularity conditions.

Assumption 2 (Regularity conditions). Suppose instrument \mathbf{z} is a centered random variable. We assume the following conditions hold:

- i. Bounded parameters: $\|\boldsymbol{\beta}\| \leq B_\beta$, $\|\boldsymbol{\Theta}\| \leq B_\Theta$.
- ii. Bounded variables: $\|\mathbf{z}\| \leq B_z$, $\|\mathbf{x}\| \leq B_x$, $|\epsilon_1| \leq B_{\epsilon_1}$, $\|\epsilon_2\| \leq B_{\epsilon_2}$.
- iii. Linear instrument: $\mathbb{E}[x_k | \mathbf{z}] = \langle \boldsymbol{\Theta}_k, \mathbf{z} \rangle$.

The boundedness condition in (ii) is required to invoke matrix Bernstein inequalities (Tropp, 2015) in the analysis. We anticipate that this condition may be relaxed to subgaussian or moment conditions by using more sophisticated matrix concentration results.

Theorem 2.1 (MSE of 2SLS estimator). *Given Assumptions 1 and 2, consider clipping operation*

$$\text{clip}_{B_\beta}(\hat{\boldsymbol{\beta}}) := \begin{cases} \hat{\boldsymbol{\beta}} & \text{if } \|\hat{\boldsymbol{\beta}}\| \leq B_\beta \\ \frac{B_\beta}{\|\hat{\boldsymbol{\beta}}\|} \hat{\boldsymbol{\beta}} & \text{if } \|\hat{\boldsymbol{\beta}}\| > B_\beta \end{cases}.$$

When

$$n \geq \max \left\{ 4c^2 B_z^4 \left(q + \log \left(\frac{4c^2 B_z^4 K}{q} \right) - \frac{3}{2} \right), \frac{qe^{\frac{3}{2}}}{K}, \frac{p^2(q+1)^2 K}{qK_0^2} \right\},$$

where $K := \frac{\lambda_{\min}(\Sigma_z)}{6B_z^2}$ and $K_0 := \frac{\lambda_{\min}(\Sigma_z)\sigma_{\min}^2(\Theta)}{2B_z^2\epsilon_2}$, the mean squared error of the 2SLS estimate is bounded by:

$$\mathbb{E} \left[\|\text{clip}_{B_\beta}(\hat{\beta}_{2SLS}) - \beta\|^2 \right] \leq \mathcal{O} \left(\frac{q}{n} \left(\frac{B_\beta^2}{K} + C^2(n)\sigma_1^2 \right) \right), \quad (3)$$

where $C(n)$ is defined in Equation (42), $\Sigma_z := \mathbb{E}[zz^\top]$, and c is an absolute constant.

Remark 2.1. We keep the slightly complicated form (3) so that the \mathcal{O} notation only hides some absolute constant multipliers that are independent of problem-related constants. Note that when n is large enough, we have $C(n) \rightarrow \frac{B_\Theta B_z}{\lambda_{\min}(\Sigma_z)\sigma_{\min}^2(\Theta)}$, so $C(n)$ is also bounded. Thus the error bound (3) decays with rate $\mathcal{O}(\frac{1}{n})$.

We note that although the consistency of the 2SLS estimator is a standard result in econometrics, most existing works focus on the asymptotic properties of the estimator. Theorem 2.1 provides the first non-asymptotic bound for estimation error $\|\hat{\beta}_{2SLS} - \beta\|^2$, under random design. The detailed proof is provided in Appendix A.1.

3 TRANSFORMERS HANDLE ENDOGENEITY

3.1 TRANSFORMER ARCHITECTURE

Denote the input matrix as $\mathbf{H} = [\mathbf{h}_1, \dots, \mathbf{h}_n] \in \mathbb{R}^{D \times n}$, where each column corresponds to one sample vector.

Definition 3.1 (Attention layer). A self-attention layer with M heads is denoted as $\text{ATTN}_\theta(\cdot)$, with parameters $\theta = \{(\mathbf{Q}_m, \mathbf{K}_m, \mathbf{V}_m)\}_{m \in [M]} \subseteq \mathbb{R}^{D \times D}$. Given input \mathbf{H} ,

$$\tilde{\mathbf{H}} = \text{ATTN}_\theta(\mathbf{H}) := \mathbf{H} + \frac{1}{n} \sum_{m=1}^M (\mathbf{V}_m \mathbf{H}) \times \sigma((\mathbf{Q}_m \mathbf{H})^\top (\mathbf{K}_m \mathbf{H})) \in \mathbb{R}^{D \times n}, \quad (4)$$

or element-wise:

$$\tilde{\mathbf{h}}_i = [\text{ATTN}_\theta(\mathbf{H})]_i := \mathbf{h}_i + \sum_{m=1}^M \frac{1}{n} \sum_{j=1}^n \sigma(\langle \mathbf{Q}_m \mathbf{h}_i, \mathbf{K}_m \mathbf{h}_j \rangle) \cdot \mathbf{V}_m \mathbf{h}_j \in \mathbb{R}^D, \quad (5)$$

where $\sigma(\cdot)$ is the ReLU function.

Definition 3.2 (MLP layer). An MLP layer is denoted as $\text{MLP}_\theta(\cdot)$, with parameters $\theta = (\mathbf{W}_1, \mathbf{W}_2) \in \mathbb{R}^{D' \times D \times D \times D'}$. Given input \mathbf{H} ,

$$\tilde{\mathbf{H}} = \text{MLP}_\theta(\mathbf{H}) := \mathbf{H} + \mathbf{W}_2 \sigma(\mathbf{W}_1 \mathbf{H}),$$

or element-wise:

$$\tilde{\mathbf{h}}_i = [\text{MLP}_\theta(\mathbf{H})]_i := \mathbf{h}_i + \mathbf{W}_2 \sigma(\mathbf{W}_1 \mathbf{h}_i).$$

Definition 3.3 (Transformer). An L-layer transformer is denoted as $\text{TF}_\theta(\cdot)$, with parameters $\theta = (\theta_{\text{ATTN}}^{(1:L)}, \theta_{\text{MLP}}^{(1:L)})$. Given input $\mathbf{H} = \mathbf{H}^{(0)}$,

$$\mathbf{H}^{(l)} = \text{MLP}_{\theta_{\text{MLP}}^{(l)}}(\text{ATTN}_{\theta_{\text{ATTN}}^{(l)}}(\mathbf{H}^{(l-1)})), \quad l = 1, \dots, L.$$

The output of this transformer is the final layer output: $\tilde{\mathbf{H}} := \mathbf{H}^{(L)} = \text{TF}_\theta(\mathbf{H}^{(0)})$.

Definition 3.4 (Looped transformer). An \bar{L} -looped transformer is a special transformer architecture, denoted as $\text{LTF}_{\bar{\theta}, \bar{L}}(\cdot)$, with parameters $\bar{\theta} = (\bar{\theta}_{\text{ATTN}}^{(1:L_0)}, \bar{\theta}_{\text{MLP}}^{(1:L_0)})$. Given input $\mathbf{H} = \mathbf{H}^{(0)}$,

$$\mathbf{H}^{(l)} = \text{TF}_{\bar{\theta}}(\mathbf{H}^{(l-1)}), \quad l = 1, \dots, \bar{L}.$$

The output of this looped transformer is the final loop output: $\tilde{\mathbf{H}} := \mathbf{H}^{(\bar{L})} = \text{LTF}_{\bar{\theta}, \bar{L}}(\mathbf{H}^{(0)})$.

Previous works (e.g., Bai et al. (2023), Zhang et al. (2024a)) have shown that transformers can perform in-context linear regression by emulating gradient descent (GD) with in-context pretraining. However, these studies have two key limitations. First, their analysis is based on single-level optimization algorithms, which is insufficient to demonstrate that transformers can efficiently learn more complex algorithms like 2SLS (Definition 2.1). Second, most ICL-related research focuses on the predictive performance of transformers, paying little attention to their ability to provide accurate coefficient estimates. We extend the current ICL framework by showing that transformers can implement a bi-level GD procedure (see Section 3.2) with looped transformer architecture (Definition 3.4), allowing them to efficiently emulate 2SLS and provide coefficient estimates that are at least as accurate as 2SLS in the presence of endogeneity (as in Equation (1)).

3.2 GRADIENT DESCENT BASED IV REGRESSION

We first introduce a gradient-based bi-level optimization procedure to obtain the 2SLS estimator in Equation (2). Given the dataset $(\mathbf{Z}, \mathbf{X}, \mathbf{Y}) = \{(z_i, \mathbf{x}_i, y_i)\}_{i=1}^n$, the objective function of IV regression can be formulated as the following bi-level optimization problem:

$$\min_{\boldsymbol{\beta}} \mathcal{L}(\boldsymbol{\beta}) = \frac{1}{n} \sum_{i=1}^n (y_i - z_i^\top \hat{\boldsymbol{\Theta}} \boldsymbol{\beta})^2, \quad \text{where} \quad \hat{\boldsymbol{\Theta}} := \arg \min_{\boldsymbol{\Theta}} \frac{1}{n} \sum_{j=1}^n (\mathbf{x}_j - z_j^\top \boldsymbol{\Theta})^2. \quad (6)$$

Consider the following gradient updates with learning rates α, η :

$$\boldsymbol{\Theta}^{(t+1)} = \boldsymbol{\Theta}^{(t)} - \eta \mathbf{Z}^\top (\mathbf{Z} \boldsymbol{\Theta}^{(t)} - \mathbf{X}), \quad (7a)$$

$$\boldsymbol{\beta}^{(t+1)} = \boldsymbol{\beta}^{(t)} - \alpha \boldsymbol{\Theta}^{(t)\top} \mathbf{Z}^\top (\mathbf{Z} \boldsymbol{\Theta}^{(t)} \boldsymbol{\beta}^{(t)} - \mathbf{Y}). \quad (7b)$$

Note that the GD-2SLS updates in Equation (7) are designed to solve Equation (6). We now show that regardless the convergence of $\boldsymbol{\Theta}^{(t)}$, the GD estimator $\boldsymbol{\beta}^{(t)}$ will always converge to the 2SLS estimator in Equation (2) with exponential rate.

Theorem 3.1 (Implementing 2SLS with gradient-based method). *Given training data $(\mathbf{Z}, \mathbf{X}, \mathbf{Y}) = \{(z_i, \mathbf{x}_i, y_i)\}_{i=1}^n$. Suppose the learning rates α, η satisfy the following conditions:*

$$0 < \alpha < \frac{2}{\sigma_{\max}^2(\mathbf{Z} \hat{\boldsymbol{\Theta}})} \quad \text{and} \quad 0 < \eta < \frac{2}{\sigma_{\max}^2(\mathbf{Z})},$$

where $\sigma_{\max}(\cdot)$ denotes the largest singular value of a matrix. Then, the GD updates in Equation (7) converge to the 2SLS estimator at an exponential rate:

$$\|\boldsymbol{\beta}^{(t)} - \hat{\boldsymbol{\beta}}_{2SLS}\| \leq \mathcal{O}(\Lambda^t),$$

where, with $\rho(\cdot)$ denoting the spectral radius of the matrix,

$$\Lambda := \max\{\gamma(\alpha), \kappa(\eta)\}, \quad \gamma(\alpha) := \rho(\mathbf{I} - \alpha \hat{\boldsymbol{\Theta}}^\top \mathbf{Z}^\top \mathbf{Z} \hat{\boldsymbol{\Theta}}), \quad \kappa(\eta) := \rho(\mathbf{I} - \eta \mathbf{Z}^\top \mathbf{Z}). \quad (8)$$

To the best of our knowledge, Theorem 3.1 provides the first theoretical result demonstrating that 2SLS can be efficiently implemented using a gradient-based method, with an exponential convergence rate. We provide the proof in Appendix B.1 and present simulation results in Appendix C.1 to examine the convergence behavior of the optimization process.

3.3 TRANSFORMERS CAN EFFICIENTLY IMPLEMENT GD-2SLS

The looped transformer architecture (Definition 3.4), as proposed by Giannou et al. (2023), introduces an efficient approach to learn iterative algorithms by cascading the same transformer block for multiple times. With the GD updates in Equation (7), we will show that there exists a looped transformer architecture that can efficiently learn the 2SLS estimator. We emphasize here that although we can implement 2SLS by sequentially attaching two separate GD iterates (each handling OLS for one stage), the overall convergence depends heavily on the convergence of the first stage estimate $\hat{\boldsymbol{\Theta}}$. Hence, significantly more number of layers are needed to ensure convergence. In addition, the advantage of looped transformer architecture cannot be fully exploited with this approach.

Theorem 3.2 (Implement a step of GD-2SLS with a transformer block). *Suppose the embedded input matrix takes the form:*

$$\mathbf{H}^{(2l)} = \begin{bmatrix} \mathbf{z}_1 & \cdots & \mathbf{z}_n & \mathbf{z}_{n+1} \\ \mathbf{x}_1 & \cdots & \mathbf{x}_n & \mathbf{x}_{n+1} \\ y_1 & \cdots & y_n & 0 \\ \Theta_{:,1}^{(l)} & \cdots & \Theta_{:,1}^{(l)} & \Theta_{:,1}^{(l)} \\ \vdots & \vdots & \vdots & \vdots \\ \Theta_{:,p}^{(l)} & \cdots & \Theta_{:,p}^{(l)} & \Theta_{:,p}^{(l)} \\ \beta^{(l)} & \cdots & \beta^{(l)} & \beta^{(l)} \\ \hat{\mathbf{x}}_1^{(l)} & \cdots & \hat{\mathbf{x}}_n^{(l)} & \hat{\mathbf{x}}_{n+1}^{(l)} \\ 1 & \cdots & 1 & 1 \\ 1 & \cdots & 1 & 0 \end{bmatrix} \in \mathbb{R}^{D \times (n+1)}. \quad (9)$$

Given $\mathbf{H}^{(2l)}$, there exists a double-layer attention-only transformer block with parameters $\theta = \theta_{ATTN}^{(2l+1:2l+2)} = \{(\mathbf{Q}_m^{(2l+1:2l+2)}, \mathbf{K}_m^{(2l+1:2l+2)}, \mathbf{V}_m^{(2l+1:2l+2)})\}_{m \in [M^{(2l+1:2l+2)}]} \subset \mathbb{R}^{D \times D}$, where the number of heads $M^{(2l+1)} = 2p$, $M^{(2l+2)} = 2(p+1)$ and embedding dimension $D = qp + 3p + q + 3$, that implements a 2SLS gradient update in Equation (7) with any given learning rates α, η :

$$\mathbf{H}^{2(l+1)} = \mathbf{TF}_{\theta_{ATTN}^{(2l+1:2l+2)}}(\mathbf{H}^{(2l)}) = \begin{bmatrix} \mathbf{z}_1 & \cdots & \mathbf{z}_n & \mathbf{z}_{n+1} \\ \mathbf{x}_1 & \cdots & \mathbf{x}_n & \mathbf{x}_{n+1} \\ y_1 & \cdots & y_n & 0 \\ \Theta_{:,1}^{(l+1)} & \cdots & \Theta_{:,1}^{(l+1)} & \Theta_{:,1}^{(l+1)} \\ \vdots & \vdots & \vdots & \vdots \\ \Theta_{:,p}^{(l+1)} & \cdots & \Theta_{:,p}^{(l+1)} & \Theta_{:,p}^{(l+1)} \\ \beta^{(l+1)} & \cdots & \beta^{(l+1)} & \beta^{(l+1)} \\ \hat{\mathbf{x}}_1^{(l+1)} & \cdots & \hat{\mathbf{x}}_n^{(l+1)} & \hat{\mathbf{x}}_{n+1}^{(l+1)} \\ 1 & \cdots & 1 & 1 \\ 1 & \cdots & 1 & 0 \end{bmatrix} \in \mathbb{R}^{D \times (n+1)}.$$

Our existence proof specifies an attention structure such that one layer updates only the first-stage estimate $\hat{\mathbf{x}}_i^{(l)}$ for all samples, followed by another layer to update the parameters $\Theta^{(l)}$ and $\beta^{(l)}$. Furthermore, as noted in the proof of Theorem 3.2 (ref. Appendix B.2), regardless of the initial values of $\Theta^{(l)}$, $\beta^{(l)}$ and $\hat{\mathbf{x}}^{(l)}$, the structures of the transformer blocks remain the same. This allows us to exploit the looped transformer architecture to significantly reduce the number of parameters and improve learning efficiency (Yang et al., 2024).

By cascading the transformer block \bar{L} times, with Theorem 3.1, one can show that transformers are able to mimic the 2SLS estimator with exponential convergence rate, as described in the following corollary.

Corollary 3.1 (Implementing GD-2SLS with looped transformer). *For any $0 < \varepsilon < 1$, given learning rates α, η , and $\Lambda \in (0, 1)$, as defined in Equation (8), there exists a transformer formulated as $\mathbf{TF}_\theta(\cdot) := \mathbf{TF}_{\theta'}(\mathbf{LTF}_{\bar{\theta}, \bar{L}}(\cdot))$, which consists of an \bar{L} -looped transformer $\mathbf{LTF}_{\bar{\theta}, \bar{L}}$ with $\bar{\theta} = \bar{\theta}_{ATTN}^{(1:2)} = \{(\bar{\mathbf{Q}}_m^{(1:2)}, \bar{\mathbf{K}}_m^{(1:2)}, \bar{\mathbf{V}}_m^{(1:2)})\}_{m \in [\bar{M}^{(1:2)}]} \subset \mathbb{R}^{D \times D}$, $\bar{L} = \lceil \mathcal{O}(\log_\Lambda(\varepsilon)) \rceil$, and a final attention layer¹ $\theta' = \theta'_{ATTN} = \{(\mathbf{Q}'_m, \mathbf{K}'_m, \mathbf{V}'_m)\}_{m \in [M']]} \subset \mathbb{R}^{D \times D}$, where $\bar{M}^{(1)} = 2p$, $\bar{M}^{(2)} = 2(p+1)$, $M' = 2$, such that given embedded input $\mathbf{H}^{(0)}$ taking the format in Equation (9), the model output satisfies:*

$$|\text{read}_y(\mathbf{TF}_\theta(\mathbf{H}^{(0)})) - \hat{\beta}_{2SLS}^\top \mathbf{x}_{n+1}| \leq B_x \varepsilon,$$

where $\text{read}_y(\cdot)$ is a function that reads the prediction \hat{y}_{n+1} from the output of the transformer.

We emphasize here that our construction differs from the implementation of Bai et al. (2023, Theorem 4) for OLS in the following aspects:

¹This layer updates the prediction $\hat{y}_{n+1} := \beta^{(\bar{L})\top} \mathbf{x}_{n+1}$, which can be constructed with 2 attention heads using the same architecture as Bai et al. (2023, Theorem 13)

- i. We apply the square loss as defined in Equation (6) to learn the 2SLS estimator, which simplifies the loss function’s sum-of-ReLU representation.
- ii. The dimension of the input embedding is $D = qp + 3p + q + 3$, where the extra dimensions store the vectorized parameters $\Theta^{(l)}, \beta^{(l)}$, and the first stage estimate $\hat{x}^{(l)}$.
- iii. We use a two-layer attention-only transformer block $\bar{\theta}$ to implement a 2SLS GD update (7), with the first layer to update the current first-stage estimate $\hat{x}^{(l)}$, and the second layer to update the parameters $\Theta^{(l)}$ and $\beta^{(l)}$.
- iv. For each transformer block, in the first layer, we equip 2 heads to update each dimension of $\hat{x}_i^{(l)} \in \mathbb{R}^p$ for all samples. In the second layer, we equip 2 heads to update each column of $\Theta^{(l)} \in \mathbb{R}^{q \times p}$ and $\beta^{(l)} \in \mathbb{R}^p$.

3.4 PRETRAINING AND EXCESS LOSS BOUND

With slightly abuse of notations, we denote the (formulated) training prompt as:

$$\mathbf{H}_k = \begin{bmatrix} \mathbf{z}_{1,k} & \cdots & \mathbf{z}_{n,k} & \mathbf{z}_{n+1,k} \\ \mathbf{x}_{1,k} & \cdots & \mathbf{x}_{n,k} & \mathbf{x}_{n+1,k} \\ y_{1,k} & \cdots & y_{n,k} & 0 \end{bmatrix} \in \mathbb{R}^{(p+q+1) \times (n+1)}, \quad k = 1, \dots, N.$$

Note that we denote each training prompt by the subscript $k = 1, \dots, N$, where N is the total number of prompts. Each training prompt consists of n labeled training samples $\{(\mathbf{z}_i, \mathbf{x}_i, y_i)\}_{i=1}^n$, and one unlabeled query sample $(\mathbf{z}_{n+1}, \mathbf{x}_{n+1})$. Our goal is to predict y_{n+1} given the context provided by the prompt.

We introduce the following ICL data generating scheme such that endogeneity occurs in the training samples, but does not extend to the query sample. Each training prompt is generated by the in-context distribution \mathcal{P} , described by Algorithm 1.

Algorithm 1 In-Context Distribution \mathcal{P}

- 1: **Parameters:** Sample size n , clipping thresholds B_z, B_x, B_y . Task parameters $\Theta, \beta, \Phi, \phi, \Sigma_z, \Sigma_u, \Sigma_\omega, \sigma_\epsilon$ from meta distribution π .
 - 2: **Output:** Training samples $\{(\mathbf{z}_i, \mathbf{x}_i, y_i)\}_{i=1}^n$, query sample $(\mathbf{z}_{n+1}, \mathbf{x}_{n+1}, y_{n+1})$.
 - 3: **for** $i = 1, \dots, n$ **do**
 - 4: **Generate:** $\mathbf{z}_i \sim \mathcal{N}(0, \Sigma_z)$, $\mathbf{u}_i \sim \mathcal{N}(0, \Sigma_u)$, $\boldsymbol{\omega}_i \sim \mathcal{N}(0, \Sigma_\omega)$, $\epsilon_i \sim \mathcal{N}(0, \sigma_\epsilon^2)$.
 - 5: **Compute:** $\mathbf{x}_i = \Theta^\top \mathbf{z}_i + \Phi^\top \mathbf{u}_i + \boldsymbol{\omega}_i$.
 - 6: **Compute:** $y_i = \beta^\top \mathbf{x}_i + \phi^\top \mathbf{u}_i + \epsilon_i$.
 - 7: **end for**
 - 8: **Generate:** $\mathbf{z}_{n+1} \sim \mathcal{N}(0, \Sigma_z)$, $\boldsymbol{\omega}_{n+1} \sim \mathcal{N}(0, \Sigma_\omega)$, $\epsilon_{n+1} \sim \mathcal{N}(0, \sigma_\epsilon^2)$.
 - 9: **Compute:** $\mathbf{x}_{n+1} = \Theta^\top \mathbf{z}_{n+1} + \boldsymbol{\omega}_{n+1}$.
 - 10: **Compute:** $y_{n+1} = \beta^\top \mathbf{x}_{n+1} + \epsilon_{n+1}$.
 - 11: **Clip:** $\mathbf{z}_i = \text{clip}_{B_z}(\mathbf{z}_i)$, $\mathbf{x}_i = \text{clip}_{B_x}(\mathbf{x}_i)$, $y_i = \text{clip}_{B_y}(y_i)$ for $i = 1, \dots, n + 1$.
-

In Algorithm 1, $\mathbf{u} \in \mathbb{R}^p$ is the source of endogenous error, $\boldsymbol{\omega} \in \mathbb{R}^p, \epsilon \in \mathbb{R}$ are the exogenous errors. Note that we have $\epsilon_{1,i} = \phi^\top \mathbf{u}_i + \epsilon_i$ and $\epsilon_{2,i} = \Phi^\top \mathbf{u}_i + \boldsymbol{\omega}_i$, corresponding to the notations in Equation (1). $\Theta \in \mathbb{R}^{q \times p}, \beta \in \mathbb{R}^p, \Phi \in \mathbb{R}^{p \times p}, \phi \in \mathbb{R}^p, \Sigma_z \in \mathbb{R}^{q \times q}, \Sigma_u \in \mathbb{R}^{p \times p}, \Sigma_\omega \in \mathbb{R}^{p \times p}, \sigma_\epsilon \in \mathbb{R}$ are task-specific parameters following meta distribution π . $\text{clip}_B(\cdot)$ is a clipping operator to bound the norm of input within radius B . We say that the in-context samples $\{(\mathbf{z}_i, \mathbf{x}_i, y_i)\}_{i=1}^{n+1}$ are drawn from the in-context distribution \mathcal{P} , and $\mathcal{P} \sim \pi$ if the task parameters $(\Theta, \beta, \Phi, \phi, \Sigma_z, \Sigma_u, \Sigma_\omega, \sigma_\epsilon)$ are sampled from π . One can check that Assumption 1 and Assumption 2(ii)(iii) are directly satisfied with the data generated from the in-context distribution \mathcal{P} .

Following the theoretical framework of (Bai et al., 2023), we define the population ICL loss²:

$$L_{\text{ICL}}(\boldsymbol{\theta}) = \mathbb{E}_\pi \mathbb{E}_{\mathcal{P}} [y_{n+1} - \text{clip}_{B_y}(\text{read}_y(\text{TF}_{\boldsymbol{\theta}}^R(\mathbf{H}^{(0)})))]^2, \quad (10)$$

²All the clipping operations are only for analytical purpose. In practice, the behavior of the trained transformer is consistent even without the clipping bounds.

where $\mathbf{H}^{(0)}$ is the embedded input as defined in Equation (9), TF_θ^R is the transformer model with parameter θ and clipping operation $\text{clip}_R(\cdot)$ applied to each layer output. For simplicity, we denote $\widetilde{\text{TF}}_\theta(\mathbf{H}) := \text{clip}_{B_y}(\text{read}_y(\text{TF}_\theta^R(\mathbf{H}^{(0)})))$.

The transformer is trained to minimize the in-context loss in Equation (10) with the following empirical loss:

$$\hat{L}_{\text{ICL}}(\theta) = \frac{1}{N} \sum_{k=1}^N (y_{n+1,k} - \widetilde{\text{TF}}_\theta(\mathbf{H}_k))^2. \quad (11)$$

We consider the following constrained optimization problem:

$$\begin{aligned} \hat{\theta} &:= \arg \min_{\theta \in \mathcal{D}_{L,M,D',B_\theta}} \hat{L}_{\text{ICL}}(\theta), \\ \mathcal{D}_{L,M,D',B_\theta} &:= \{\theta = (\theta_{\text{Attn}}^{(1:L)}, \theta_{\text{MLP}}^{(1:L)}) : \max_{l \in [L]} M^{(l)} \leq M, \max_{l \in [L]} D^{(l)} \leq D', \|\theta\| \leq B_\theta\}, \end{aligned} \quad (12)$$

where $\|\theta\| := \max_{l \in [L]} \{\max_{m \in [M]} \{\|Q_m^{(l)}\|, \|K_m^{(l)}\|\} + \sum_{m=1}^M \|V_m^{(l)}\| + \|W_1^{(l)}\| + \|W_2^{(l)}\|\}$.

We now establish excess loss bound for the trained transformer model.

Theorem 3.3 (Excess loss bound for in-context pretrained transformer). *Suppose condition (i) in Assumption 2 holds and the meta distribution π satisfies the following conditions:*

$$\mathbb{E}_\pi [\phi^\top \Sigma_u \phi + \sigma_\epsilon^2] \leq \tilde{\sigma}^2 \text{ and } \mathbb{E}_\pi [\sigma_\epsilon^2] \leq \tilde{\sigma}_\epsilon^2. \quad (13)$$

Let the in-context distribution $\mathcal{P} \sim \pi$ such that the samples $(z_i, \mathbf{x}_i, y_i)_{i=1}^{n+1}$ are drawn independently from \mathcal{P} (ref. Algorithm 1). With training prompts $\mathbf{H}_k, k = 1, \dots, N$, under ICL loss (10), the trained transformer (12) with $L = 2\bar{L} + 1, M = 2(p + 1), D = qp + 3p + q + 3, D' = 0$ (attention-only) achieves the following excess loss with probability at least $1 - \zeta$:

$$\begin{aligned} L_{\text{ICL}}(\hat{\theta}) - \mathbb{E}_\pi \mathbb{E}_{\mathcal{P}} [(y_{n+1} - \langle \beta, \mathbf{x}_{n+1} \rangle)^2] &\leq \mathcal{O} \left((\Lambda^*)^{\bar{L}} \left(B_x^2 \sqrt{\frac{q}{n} \left(\frac{B_\beta^2}{K} + C^2(n) \tilde{\sigma}^2 \right)} + B_x \tilde{\sigma}_\epsilon \right) \right. \\ &\quad \left. + B_x^2 \left(\frac{q}{n} \left(\frac{B_\beta^2}{K} + C^2(n) \tilde{\sigma}^2 \right) + \mu_{\Lambda,2}^* \right) + B_y^2 \sqrt{\frac{L^2 M D^2 \log(2 + \max\{B_\theta, R, B_y\}) + \log(1/\zeta)}{N}} \right), \end{aligned}$$

where $\Lambda^* := \min_{\alpha, \eta} \mathbb{E}_\pi \mathbb{E}_{\mathcal{P}} [\Lambda | \mathbf{H}, \alpha, \eta] < 1$, and $\mu_{\Lambda,2}^* := \mathbb{E}_\pi \mathbb{E}_{\mathcal{P}} [\Lambda^{2\bar{L}} | \mathbf{H}, \alpha^*, \eta^*]$ is close to 0.

In practical training, the number of prompts N is usually large enough such that the last term of the above bound is negligible. Thus, given a meta distribution π , the excess loss is dominated by two factors: (i) number of attention layers, and (ii) number of in-context samples. The proof of Theorem 3.3 is provided in Appendix B.3.

3.5 EXTRACTING THE REGRESSION COEFFICIENTS

The primary goal of IV regression is to estimate the causal effect, i.e. the coefficient β under the stated endogeneity in Equation (1). For 2SLS, the estimated causal effect is given by the coefficients of the endogenous variable in the second stage regression (2). For transformer models, we propose a straightforward method to extract these estimated coefficients by differentiating the output with respect to each dimension of the endogenous variable. The specific approach is summarized in Algorithm 2. We observe that the choice of Δ within a reasonable range does not significantly affect the estimation of the coefficients. In practice, usually a slightly larger Δ (for example $\Delta = 5$) can lead to a more stable estimation, which is possibly due to the elimination of rounding errors during computation.

4 EXPERIMENTS

4.1 EXPERIMENT SETUP

We conduct a simulation study to evaluate the performance of the ICL-pretrained transformer model in handling endogeneity. We set the maximum input sample size to 51 ($n = 50$ training samples and

Algorithm 2 Extracting the regression coefficients

-
- 1: **Input:** Trained transformer model $\text{TF}_{\hat{\theta}}$, input matrix \mathbf{H} , perturbation Δ .
 - 2: **Output:** Estimated coefficient $\hat{\beta}$.
 - 3: **Procedure:**
 - 4: Compute the output of the transformer model: $\hat{\mathbf{Y}} = \widetilde{\text{TF}}_{\hat{\theta}}(\mathbf{H})$.
 - 5: **for** each dimension $k = 1, \dots, p$ **do**
 - 6: Copy $\mathbf{H}_{\Delta(k)} = \mathbf{H}$. Set the k -th dimension of \mathbf{x}_{n+1} to be $(\mathbf{x}_{n+1})_k + \Delta$ for $\mathbf{H}_{\Delta(k)}$.
 - 7: Compute the new output value: $\hat{\mathbf{Y}}_{\Delta(k)} = \widetilde{\text{TF}}_{\hat{\theta}}(\mathbf{H}_{\Delta(k)})$.
 - 8: Compute the estimated coefficient: $\hat{\beta}_k = \frac{\hat{\mathbf{Y}}_{\Delta(k)} - \hat{\mathbf{Y}}}{\Delta}$.
 - 9: **end for**
-

one query sample), the dimension of endogenous variable $p = 5$, and the dimension of instrument $q = 10$. The training prompts are generated using Algorithm 1, with task parameters $\Theta, \beta, \Phi, \phi$ sampled from standard Gaussian distribution, and the covariance matrices $\Sigma_z, \Sigma_u, \Sigma_\omega$ set to be identity matrices. The noise level σ_ϵ is set to 1. We ignore all the clipping bounds in the experiment ($B_\beta, B_\Theta, B_z, B_x, B_y, B_\theta, R$ set to infinity).

The backbone of the transformer block is initialized using GPT-2 settings, with 12 attention heads ($M = 12$), 80-dimensional embedding space ($D = 80$) and 2 layers ($L_0 = 2$), following the theoretical guidelines in Theorem 3.2. We employ the looped transformer architecture, consisting of 10 identical cascading transformer blocks. The transformer model is trained under the ICL loss (11) with a batch size of $N = 64$, over a total of 300,000 training steps.

We evaluate the trained transformer model on test prompts that are not included during training. As benchmarks, we compare the transformer’s performance against the 2SLS and the OLS estimators, which are obtained by directly fitting the training samples $\{(z_i, \mathbf{x}_i, y_i)\}_{i=1}^n$ within the text prompts. In contrast, the same trained transformer model is used without any parameter adjustments for each task. We compare the performance of these models from two aspects: the in-context prediction error (ICPE) on the query sample y_{n+1} , and the mean squared error (MSE) on the coefficient β .

4.2 RESULTS

We first investigate the performance of the trained transformer model over endogeneity tasks with varying training sample sizes from 20 to 50. The results are shown in Figure 1a. Under endogeneity, our transformer model achieves similar performance to that of the 2SLS estimator, with only small gaps in ICPE and MSE, both outperforming the OLS estimator.

Next, we examine the performance of the trained transformer model in handling varying levels of IV strength. The strength of an instrument is measured by the correlation between the IV and the endogenous variable. To vary the IV strength, we generate prompts with z_i and \mathbf{x}_i following different correlation levels. Specifically, in Algorithm 1, we adjust the IV strength by multiplying Θ by a factor $r \in (0, 2)$ when generating test prompts. The results are shown in Figure 1b.

Interestingly, the trained transformer model outperforms the 2SLS estimator in handling weaker IVs (when IV strength < 0.5). This suggests that, beyond merely mimicking 2SLS, the ICL training process may equip the transformer model with a more advanced mechanism for handling endogeneity with weak IVs than the 2SLS estimator. At the same time, when the IV is strong, the transformer model maintains performance comparable to that of the 2SLS estimator.

This finding motivates us to further examine the performance of the trained transformer model in non-standard endogeneity tasks. We consider two scenarios: (a) the IV has a quadratic effect on the endogenous variable, i.e. $\mathbf{x}_{i,k} = \Theta_k^\top z_{i,k}^2 + \text{error}_{i,k}$ in Algorithm 1, and (b) the dimension of IV is not sufficient to identify the endogenous variable³, where we set $q = 3$ (by zeroing out the remaining dimensions of z in test prompts) and $p = 5$.

We evaluate the same trained transformer model as before, with results presented in Figure 2a and Figure 2b, respectively. Once again, the trained transformer model consistently outperforms both

³For 2SLS estimate, the actual computation uses pseudoinverse to handle rank deficiency.

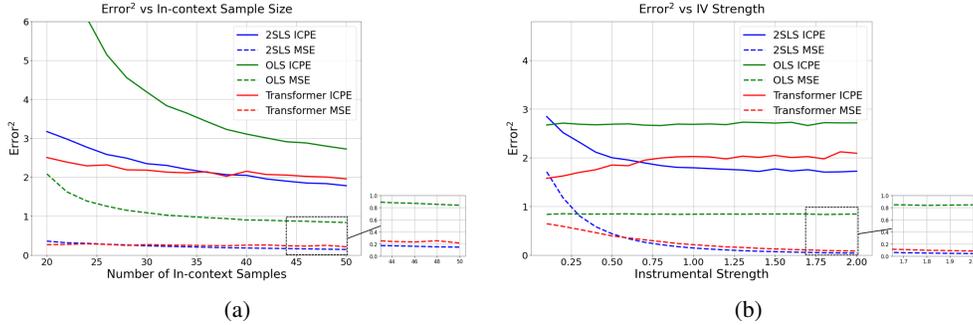


Figure 1: The ICL performance of the trained transformer model in endogeneity tasks. We compare in-context prediction error (ICPE) and coefficient MSE versus (a) the number of in-context samples; (b) the IV strength. The curves are averaged over 500 simulations.

2SLS and OLS estimators in handling these non-standard endogeneity tasks. All these results suggest that the trained transformer can be generalized effectively to a broader range of endogeneity tasks while still providing reliable in-context predictions and coefficient estimates. To further illustrate this capability, we also examine other cases including multicollinearity, complex non-linear IV, and varying endogeneity strengths, see Appendix C.3,C.4,C.5. We suspect that, in our pretraining scheme, although the 2SLS estimator already achieves small excess loss, a gap remains between the 2SLS estimator and the optimal predictor that the transformer model successfully bridges. Finally, we conclude that through ICL training, the transformer model performs at least as well as 2SLS and appears to be a promising tool for handling endogeneity in difficult scenarios.

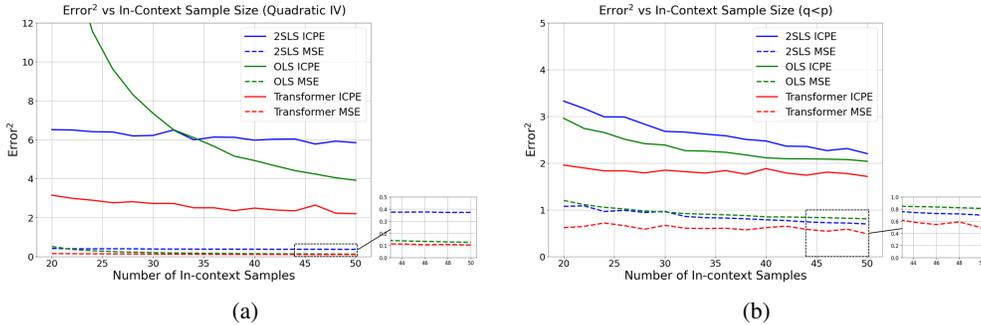


Figure 2: The ICL performance of the trained transformer model in non-standard endogeneity tasks: (a) The IV has quadratic effect on the endogenous variable; (b) The dimension of IV is not sufficient to identify the endogenous variable. The curves are averaged over 500 simulations.

5 CONCLUSION

This paper presents a novel perspective on the transformer model in its ability to handle endogeneity in in-context linear regression. We have theoretically shown that the transformer model exists an intrinsic structure that enables it to learn the 2SLS algorithm through an efficient GD procedure. We have further provided a theoretical guarantee that the trained transformer model can achieve a small excess loss over the optimal loss, under our proposed ICL training scheme. Our simulation study demonstrates that the trained transformer model can achieve comparable performance to the 2SLS estimator in handling standard endogeneity tasks. Furthermore, our investigation illustrates that it exhibits significantly better performances in handling complex scenarios such as weak instruments, non-linear IV, and underdetermined IV problems, compared to the 2SLS estimator. These results suggest that the ICL pre-trained transformer model is a promising tool for making reliable in-context predictions and coefficient estimates under endogeneity, especially when dealing with non-standard IV problems.

ACKNOWLEDGEMENTS

The work of H. Liang and L. Lai was supported by National Science Foundation (NSF) under grants CCF-2112504 and CCF-2232907. The work of K. Balasubramanian was supported by NSF under grants DMS-2413426 and DMS-2053918.

REFERENCES

- Kwangjun Ahn, Xiang Cheng, Hadi Daneshmand, and Suvrit Sra. Transformers learn to implement preconditioned gradient descent for in-context learning. In *Proceedings of Neural Information Processing Systems*, New Orleans, LA, USA, December 2023. URL <https://arxiv.org/abs/2306.00297>.
- Ekin Akyürek, Dale Schuurmans, Jacob Andreas, Tengyu Ma, and Denny Zhou. What learning algorithm is in-context learning? investigations with linear models. In *Proceedings of International Conference on Learning Representations*, Kigali, Rwanda, May 2023. URL <https://openreview.net/forum?id=0g0X4H8yN4I>.
- Joshua D. Angrist and William N. Evans. Children and their parents’ labor supply: Evidence from exogenous variation in family size. *The American Economic Review*, 88(3):450–477, 1998. ISSN 00028282. URL <http://www.jstor.org/stable/116844>.
- Joshua D. Angrist and Alan B. Krueger. Instrumental variables and the search for identification: From supply and demand to natural experiments. *Journal of Economic Perspectives*, 15(4):69–85, 2001. URL <https://www.aeaweb.org/articles?id=10.1257/jep.15.4.69>.
- Joshua D Angrist and Jörn-Steffen Pischke. *Mostly harmless econometrics: An empiricist’s companion*. Princeton University Press, 2009. ISBN 9780691120355. URL <https://press.princeton.edu/books/hardcover/9780691120355/mostly-harmless-econometrics>.
- Yu Bai, Fan Chen, Huan Wang, Caiming Xiong, and Song Mei. Transformers as statisticians: Provable in-context learning with in-context algorithm selection. In *Proceedings of Neural Information Processing Systems*, New Orleans, LA, USA, December 2023. URL https://proceedings.neurips.cc/paper_files/paper/2023/file/b2e63e36c57e153b9015fece2352a9f9-Paper-Conference.pdf.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. In *Proceedings of Neural Information Processing Systems*, Virtual, December 2020. URL <https://arxiv.org/abs/2005.14165>.
- Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *Proceedings of European Conference on Computer Vision*, Virtual, August 2020. URL <https://arxiv.org/abs/2005.12872>.
- Xuxing Chen, Abhishek Roy, Yifan Hu, and Krishnakumar Balasubramanian. Stochastic optimization algorithms for instrumental variable regression with streaming data. In *Proceedings of Neural Information Processing Systems*, Vancouver, Canada, December 2024. URL <https://arxiv.org/abs/2405.19463>.
- Xiang Cheng, Yuxin Chen, and Suvrit Sra. Transformers implement functional gradient descent to learn non-linear functions in context. In *Proceedings of International Conference on Machine Learning*, Honolulu, HI, USA, July 2024. URL <https://proceedings.mlr.press/v235/cheng24a.html>.

- Riccardo Della Vecchia and Debabrota Basu. Stochastic online instrumental variable regression: Regrets for endogeneity and bandit feedback. *arXiv preprint arXiv:2302.09357*, 2023. URL <https://arxiv.org/abs/2302.09357>.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the North American Chapter of the Association for Computational Linguistics*, Minneapolis, MN, USA, June 2019. URL <https://api.semanticscholar.org/CorpusID:52967399>.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *Proceedings of International Conference on Learning Representations*, Virtual, May 2021. URL <https://openreview.net/forum?id=YicbFdNTTy>.
- Benjamin L Edelman, Ezra Edelman, Surbhi Goel, Eran Malach, and Nikolaos Tsilivis. The evolution of statistical induction heads: In-context learning markov chains. In *Proceedings of Neural Information Processing Systems*, Vancouver, Canada, December 2024. URL <https://arxiv.org/abs/2402.11004>.
- Fabian Falck, Ziyu Wang, and Chris Holmes. Is in-context learning in large language models bayesian? a martingale perspective. *arXiv preprint arXiv:2406.00793*, 2024. URL <https://arxiv.org/abs/2406.00793>.
- Yuri Fonseca, Caio Peixoto, and Yuri Saporito. Nonparametric instrumental variable regression through stochastic approximate gradients. In *Proceedings of Neural Information Processing Systems*, Vancouver, Canada, December 2024. URL <https://arxiv.org/abs/2402.05639>.
- Deqing Fu, Tian-Qi Chen, Robin Jia, and Vatsal Sharan. Transformers learn higher-order optimization methods for in-context learning: A study with linear models. *arXiv preprint arXiv:2310.17086*, 2023. URL <https://arxiv.org/abs/2310.17086>.
- Yihang Gao, Chuanyang Zheng, Enze Xie, Han Shi, Tianyang Hu, Yu Li, Michael K. Ng, Zhenguo Li, and Zhaoqiang Liu. On the expressive power of a variant of the looped transformer. *arXiv preprint arXiv:2402.13572*, 2024. URL <https://arxiv.org/abs/2402.13572>.
- Shivam Garg, Dimitris Tsipras, Percy S Liang, and Gregory Valiant. What can transformers learn in-context? a case study of simple function classes. In *Proceedings of Neural Information Processing Systems*, New Orleans, LA, USA, December 2022. URL <https://arxiv.org/abs/2208.01066>.
- Angeliki Giannou, Shashank Rajput, Jy-Yong Sohn, Kangwook Lee, Jason D. Lee, and Dimitris Papailiopoulos. Looped transformers as programmable computers. In *Proceedings of International Conference on Machine Learning*, Honolulu, HI, USA, July 2023. URL <https://proceedings.mlr.press/v202/giannou23a.html>.
- Angeliki Giannou, Liu Yang, Tianhao Wang, Dimitris Papailiopoulos, and Jason D. Lee. How well can transformers emulate in-context newton’s method? *arXiv preprint arXiv:2403.03183*, 2024. URL <https://arxiv.org/abs/2403.03183>.
- Gautam Goel and Peter Bartlett. Can a transformer represent a kalman filter? In *Proceedings of Annual Learning for Dynamics & Control Conference*, Oxford, UK, July 2024. URL <https://proceedings.mlr.press/v242/goel24a.html>.
- William H. Greene. *Econometric Analysis*. Pearson, 8th edition, 2018. ISBN 978-0-13-446136-6. URL <https://pages.stern.nyu.edu/~wgreene/Text/econometricanalysis.htm>.
- Jerry Hausman. Mismeasured variables in econometric analysis: problems from the right and problems from the left. *Journal of Economic perspectives*, 15(4):57–67, 2001. URL <https://www.aeaweb.org/articles?id=10.1257/jep.15.4.57>.

- Daniel Hsu, Sham M. Kakade, and Tong Zhang. Random design analysis of ridge regression. *Foundations of Computational Mathematics*, 14(3):569–600, 2014. URL <https://link.springer.com/article/10.1007/s10208-014-9192-1>.
- Yu Huang, Yuan Cheng, and Yingbin Liang. In-context convergence of transformers. In *Proceedings of International Conference on Machine Learning*, Honolulu, HI, USA, July 2024. URL <https://proceedings.mlr.press/v235/huang24d.html>.
- Yanhao Jin, Krishnakumar Balasubramanian, and Debashis Paul. Meta-learning with generalized ridge regression: High-dimensional asymptotics, optimality and hyper-covariance estimation. *arXiv preprint arXiv:2403.19720*, 2024. URL <https://arxiv.org/abs/2403.19720>.
- Yanhao Jin, Krishnakumar Balasubramanian, and Lifeng Lai. In-context learning for mixture of linear regressions: Existence, generalization and training dynamics. *arXiv preprint arXiv:2410.14183*, 2025. URL <https://arxiv.org/abs/2410.14183>.
- Kuang-Huei Lee, Ofir Nachum, Mengjiao (Sherry) Yang, Lisa Lee, Daniel Freeman, Sergio Guadarrama, Ian Fischer, Winnie Xu, Eric Jang, Henryk Michalewski, and Igor Mordatch. Multi-game decision transformers. In *Proceedings of Neural Information Processing Systems*, New Orleans, LA, USA, November 2022. URL <https://arxiv.org/abs/2211.15196>.
- Yingcong Li, Muhammed Emrullah Ildiz, Dimitris Papailiopoulos, and Samet Oymak. Transformers as algorithms: Generalization and stability in in-context learning. In *Proceedings of International Conference on Machine Learning*, Honolulu, HI, USA, July 2023. URL <https://proceedings.mlr.press/v202/li23l.html>.
- Ashok Vardhan Makkuva, Marco Bondaschi, Chanakya Ekbote, Adway Girish, Alliot Nagle, Hyeji Kim, and Michael Gastpar. Local to global: Learning dynamics and effect of initialization for transformers. In *ICML 2024 Workshop on Theoretical Foundations of Foundation Models*, Vienna, Austria, July 2024a. URL <https://openreview.net/forum?id=OYoCJPwbfc>.
- Ashok Vardhan Makkuva, Marco Bondaschi, Alliot Nagle, Adway Girish, Hyeji Kim, Martin Jaggi, and Michael Gastpar. Attention with markov: A curious case of single-layer transformers. In *ICML 2024 Workshop on Mechanistic Interpretability*, Vienna, Austria, July 2024b. URL <https://openreview.net/forum?id=xi6lie0SUr>.
- Krikamol Muandet, Arash Mehrjou, Si Kai Lee, and Anant Raj. Dual instrumental variable regression. In *Proceedings of Neural Information Processing Systems*, Vancouver, Canada, December 2020. URL <https://proceedings.neurips.cc/paper/2020/hash/1c383cd30b7c298ab50293adfecb7b18-Abstract.html>.
- Eshaan Nichani, Alex Damian, and Jason D. Lee. How transformers learn causal structure with gradient descent. In *Proceedings of International Conference on Machine Learning*, Honolulu, HI, USA, July 2024. URL <https://arxiv.org/abs/2402.14735>.
- Emilio Parisotto, H. Francis Song, Jack W. Rae, Razvan Pascanu, Caglar Gulcehre, Siddhant M. Jayakumar, Max Jaderberg, Raphael Lopez Kaufman, Aidan Clark, Seb Noury, Matthew M. Botvinick, Nicolas Heess, and Raia Hadsell. Stabilizing transformers for reinforcement learning. In *Proceedings of International Conference on Machine Learning*, Virtual, July 2020. URL <https://proceedings.mlr.press/v119/parisotto20a.html>.
- Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners. In *OpenAI Technical Report*, 2019. URL <https://api.semanticscholar.org/CorpusID:160025533>.
- Nived Rajaraman, Marco Bondaschi, Kannan Ramchandran, Michael Gastpar, and Ashok Vardhan Makkuva. Transformers on markov data: Constant depth suffices. *arXiv preprint arXiv:2407.17686*, 2024. URL <https://arxiv.org/abs/2407.17686>.
- Rahul Singh, Maneesh Sahani, and Arthur Gretton. Kernel instrumental variable regression. In *Proceedings of Neural Information Processing Systems*, Vancouver, Canada, December 2019. URL <https://papers.nips.cc/paper/8708-kernel-instrumental-variable-regression>.

- J.H. Stock and M.W. Watson. *Introduction to Econometrics*. Addison-Wesley, 3rd edition, 2011. ISBN 9780138009007. URL <https://stock.scholars.harvard.edu/publications/introduction-econometrics-0>.
- Joel A. Tropp. *An Introduction to Matrix Concentration Inequalities*. Now Publishers Inc., 2015. ISBN 978-1-60198-838-6. URL <https://doi.org/10.1561/22000000048>.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Proceedings of Neural Information Processing Systems*, Long Beach, CA, USA, December 2017. URL https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf.
- Roman Vershynin. *High-Dimensional Probability: An Introduction with Applications in Data Science*. Cambridge University Press, 2018. ISBN 978-1-108-41519-4. URL <https://www.cambridge.org/core/books/highdimensional-probability/797C466DA29743D2C8213493BD2D2102>.
- Johannes Von Oswald, Eyvind Niklasson, Ettore Randazzo, João Sacramento, Alexander Mordvintsev, Andrey Zhmoginov, and Max Vladymyrov. Transformers learn in-context by gradient descent. In *Proceedings of International Conference on Machine Learning*, Honolulu, HI, USA, July 2023. URL <https://proceedings.mlr.press/v202/von-oswald23a.html>.
- Charles F. Westoff and Robert Parke. *Demographic and social aspects of population growth*. Commission on Population Growth and the American Future, 1972. URL <https://catalog.hathitrust.org/Record/000008850>.
- J.M. Wooldridge. *Introductory Econometrics: A Modern Approach*. Cengage Learning, 2015. ISBN 9781473754393. URL <https://books.google.com/books?id=HveHAQAACAAJ>.
- Liu Yang, Kangwook Lee, Robert Nowak, and Dimitris Papailiopoulos. Looped transformers are better at learning learning algorithms. In *Proceedings of International Conference on Learning Representations*, Vienna, Austria, May 2024. URL <https://openreview.net/forum?id=HHbRxoDTxE>.
- Naimeng Ye, Hanming Yang, Andrew Siah, and Hongseok Namkoong. Pre-training and in-context learning is bayesian inference a la de finetti. In *ICLR 2024 Workshop on Mathematical and Empirical Understanding of Foundation Models*, Vienna, Austria, May 2024. URL <https://openreview.net/forum?id=ttupfosvgx>.
- Ruiqi Zhang, Spencer Frei, and Peter L Bartlett. Trained transformers learn linear models in-context. *Journal of Machine Learning Research*, 25(49):1–55, 2024a. URL <https://www.jmlr.org/papers/v25/23-1042.html>.
- Ruiqi Zhang, Jingfeng Wu, and Peter L Bartlett. In-context learning of a linear transformer block: benefits of the mlp component and one-step gd initialization. In *Proceedings of Neural Information Processing Systems*, New Orleans, LA, USA, December 2024b. URL <https://arxiv.org/abs/2402.14951>.

A PROOFS FOR SECTION 2

A.1 PROOF OF THEOREM 2.1

We first introduce the following lemmas that are used in the proof of Theorem 2.1.

Lemma A.1 (Bernstein Inequality, from Theorem 6.1.1 in Tropp (2015)). Let $\mathbf{S}_1, \dots, \mathbf{S}_n$ be independent, centered random matrices with common dimension $d_1 \times d_2$, and assume that each one is almost surely bounded:

$$\mathbb{E}[\mathbf{S}_i] = \mathbf{0}, \mathbb{P}(\|\mathbf{S}_i\| \leq b) = 1, \quad \forall i = 1, \dots, n.$$

With the sum:

$$\mathbf{\Omega} = \sum_{i=1}^n \mathbf{S}_i,$$

and the matrix variance statistic of the sum:

$$\nu(\mathbf{\Omega}) := \max \{ \|\mathbb{E}(\mathbf{\Omega}\mathbf{\Omega}^\top)\|, \|\mathbb{E}(\mathbf{\Omega}^\top\mathbf{\Omega})\| \},$$

then the following inequality holds:

$$\mathbb{P}\{\|\mathbf{\Omega}\| \geq \varepsilon\} \leq (d_1 + d_2) \cdot \exp\left(\frac{-\varepsilon^2/2}{\nu(\mathbf{\Omega}) + b\varepsilon/3}\right) \text{ for any } \varepsilon \geq 0.$$

Lemma A.2 (Inverse Convergence, adapted from Lemma 2.1 in Jin et al. (2024)). Suppose we have a random invertible matrix $\mathbf{\Omega}$ and invertible matrix sequence $\{\hat{\mathbf{\Omega}}^{(n)}\}$ such that $\hat{\mathbf{\Omega}}^{(n)} \xrightarrow{\mathbb{P}} \mathbf{\Omega}$. If there exists a constant $\tilde{\lambda} > 0$ such that $\sigma_{\min}(\hat{\mathbf{\Omega}}) \geq \tilde{\lambda}$ almost surely, then it holds that:

$$(\hat{\mathbf{\Omega}}^{(n)})^{-1} \xrightarrow{\mathbb{P}} \mathbf{\Omega}^{-1}.$$

Further, given convergence rate

$$\mathbb{P}\left\{\left\|\hat{\mathbf{\Omega}}^{(n)} - \mathbf{\Omega}\right\| \geq \varepsilon\right\} \leq \xi(n, \varepsilon),$$

then:

$$\mathbb{P}\left\{\left\|(\hat{\mathbf{\Omega}}^{(n)})^{-1} - \mathbf{\Omega}^{-1}\right\| \geq \varepsilon\right\} \leq \xi(n, \tilde{\lambda}^2\varepsilon).$$

Proof. We have the following decomposition:

$$(\hat{\mathbf{\Omega}}^{(n)})^{-1} - \mathbf{\Omega}^{-1} = (\hat{\mathbf{\Omega}}^{(n)})^{-1}(\mathbf{\Omega} - \hat{\mathbf{\Omega}}^{(n)})\mathbf{\Omega}^{-1}.$$

It follows that:

$$\begin{aligned} \left\|(\hat{\mathbf{\Omega}}^{(n)})^{-1} - \mathbf{\Omega}^{-1}\right\| &\leq \left\|(\hat{\mathbf{\Omega}}^{(n)})^{-1}\right\| \left\|\mathbf{\Omega} - \hat{\mathbf{\Omega}}^{(n)}\right\| \left\|\mathbf{\Omega}^{-1}\right\| \\ &\leq \frac{1}{\tilde{\lambda}^2} \left\|\mathbf{\Omega} - \hat{\mathbf{\Omega}}^{(n)}\right\|. \end{aligned}$$

Then

$$\begin{aligned} \mathbb{P}\left\{\left\|(\hat{\mathbf{\Omega}}^{(n)})^{-1} - \mathbf{\Omega}^{-1}\right\| \geq \varepsilon\right\} &\leq \mathbb{P}\left\{\frac{1}{\tilde{\lambda}^2} \left\|\mathbf{\Omega} - \hat{\mathbf{\Omega}}^{(n)}\right\| \geq \varepsilon\right\} \\ &\leq \xi(n, \tilde{\lambda}^2\varepsilon). \end{aligned}$$

□

Lemma A.3 (Product Convergence). Let $\{\hat{\mathbf{\Omega}}_1^{(n)}\}, \{\hat{\mathbf{\Omega}}_2^{(n)}\}, \dots, \{\hat{\mathbf{\Omega}}_K^{(n)}\}$ be K sequences of matrices such that $\hat{\mathbf{\Omega}}_1^{(n)} \xrightarrow{\mathbb{P}} \mathbf{\Omega}_1, \hat{\mathbf{\Omega}}_2^{(n)} \xrightarrow{\mathbb{P}} \mathbf{\Omega}_2, \dots, \hat{\mathbf{\Omega}}_K^{(n)} \xrightarrow{\mathbb{P}} \mathbf{\Omega}_K$, where each $\|\hat{\mathbf{\Omega}}_k^{(n)}\|$ is almost surely bounded for every $k = 1, \dots, K$. If the dimensions match, then it holds that:

$$\hat{\mathbf{\Omega}}_1^{(n)} \hat{\mathbf{\Omega}}_2^{(n)} \dots \hat{\mathbf{\Omega}}_K^{(n)} \xrightarrow{\mathbb{P}} \mathbf{\Omega}_1 \mathbf{\Omega}_2 \dots \mathbf{\Omega}_K.$$

Further, given convergence rates:

$$\begin{aligned} \mathbb{P} \left\{ \left\| \hat{\Omega}_1^{(n)} - \Omega_1 \right\| \geq \varepsilon \right\} &\leq \xi_1(n, \varepsilon), \\ \mathbb{P} \left\{ \left\| \hat{\Omega}_2^{(n)} - \Omega_2 \right\| \geq \varepsilon \right\} &\leq \xi_2(n, \varepsilon), \\ &\vdots \\ \mathbb{P} \left\{ \left\| \hat{\Omega}_K^{(n)} - \Omega_K \right\| \geq \varepsilon \right\} &\leq \xi_K(n, \varepsilon), \end{aligned}$$

then it holds that:

$$\mathbb{P} \left\{ \left\| \hat{\Omega}_1^{(n)} \hat{\Omega}_2^{(n)} \cdots \hat{\Omega}_K^{(n)} - \Omega_1 \Omega_2 \cdots \Omega_K \right\| \geq \varepsilon \right\} \leq \sum_{i=1}^K \xi_i \left(n, \frac{\varepsilon}{K \prod_{k \neq i} M_k} \right), \quad (14)$$

where M_k is an upper bound such that $\|\hat{\Omega}_k^{(n)}\| \leq M_k$ almost surely, $\forall k = 1, \dots, K$.

Proof. We begin by showing the case of $K = 2$. By the triangle inequality, we have:

$$\begin{aligned} \left\| \hat{\Omega}_1^{(n)} \hat{\Omega}_2^{(n)} - \Omega_1 \Omega_2 \right\| &\leq \left\| \hat{\Omega}_1^{(n)} \hat{\Omega}_2^{(n)} - \Omega_1 \hat{\Omega}_2^{(n)} \right\| + \left\| \Omega_1 \hat{\Omega}_2^{(n)} - \Omega_1 \Omega_2 \right\| \\ &\leq \left\| \hat{\Omega}_2^{(n)} \right\| \left\| \hat{\Omega}_1^{(n)} - \Omega_1 \right\| + \left\| \hat{\Omega}_2^{(n)} - \Omega_2 \right\| \left\| \Omega_1 \right\| \\ &\leq M_2 \left\| \hat{\Omega}_1^{(n)} - \Omega_1 \right\| + M_1 \left\| \hat{\Omega}_2^{(n)} - \Omega_2 \right\|. \end{aligned}$$

Using the union bound, we have:

$$\begin{aligned} &\mathbb{P} \left\{ \left\| \hat{\Omega}_1^{(n)} \hat{\Omega}_2^{(n)} - \Omega_1 \Omega_2 \right\| \geq \varepsilon \right\} \\ &\leq \mathbb{P} \left\{ M_2 \left\| \hat{\Omega}_1^{(n)} - \Omega_1 \right\| + M_1 \left\| \hat{\Omega}_2^{(n)} - \Omega_2 \right\| \geq \varepsilon \right\} \\ &\leq \mathbb{P} \left\{ M_2 \left\| \hat{\Omega}_1^{(n)} - \Omega_1 \right\| \geq \varepsilon/2 \right\} + \mathbb{P} \left\{ M_1 \left\| \hat{\Omega}_2^{(n)} - \Omega_2 \right\| \geq \varepsilon/2 \right\} \\ &\leq \xi_1 \left(n, \frac{\varepsilon}{2M_2} \right) + \xi_2 \left(n, \frac{\varepsilon}{2M_1} \right). \end{aligned}$$

For any $K > 2$, suppose the statement (14) holds for $k = 2, \dots, K-1$. Observe that:

$$\begin{aligned} &\left\| \hat{\Omega}_1^{(n)} \hat{\Omega}_2^{(n)} \cdots \hat{\Omega}_K^{(n)} - \Omega_1 \Omega_2 \cdots \Omega_K \right\| \\ &\leq \left\| \hat{\Omega}_1^{(n)} \hat{\Omega}_2^{(n)} \cdots \hat{\Omega}_K^{(n)} - \Omega_1 \Omega_2 \cdots \Omega_{K-1} \hat{\Omega}_K^{(n)} \right\| + \left\| \Omega_1 \Omega_2 \cdots \Omega_{K-1} \hat{\Omega}_K^{(n)} - \Omega_1 \Omega_2 \cdots \Omega_K \right\| \\ &\leq M_K \left\| \hat{\Omega}_1^{(n)} \hat{\Omega}_2^{(n)} \cdots \hat{\Omega}_{K-1}^{(n)} - \Omega_1 \Omega_2 \cdots \Omega_{K-1} \right\| + \prod_{k=1}^{K-1} M_k \left\| \hat{\Omega}_K^{(n)} - \Omega_K \right\|. \end{aligned} \quad (15)$$

Then it follows that:

$$\begin{aligned}
& \mathbb{P} \left\{ \left\| \hat{\Omega}_1^{(n)} \hat{\Omega}_2^{(n)} \cdots \hat{\Omega}_K^{(n)} - \Omega_1 \Omega_2 \cdots \Omega_K \right\| \geq \varepsilon \right\} \\
& \leq \mathbb{P} \left\{ M_K \left\| \hat{\Omega}_1^{(n)} \hat{\Omega}_2^{(n)} \cdots \hat{\Omega}_{K-1}^{(n)} - \Omega_1 \Omega_2 \cdots \Omega_{K-1} \right\| + \prod_{k=1}^{K-1} M_k \left\| \hat{\Omega}_K^{(n)} - \Omega_K \right\| \geq \varepsilon \right\} \\
& \leq \mathbb{P} \left\{ M_K \left\| \hat{\Omega}_1^{(n)} \hat{\Omega}_2^{(n)} \cdots \hat{\Omega}_{K-1}^{(n)} - \Omega_1 \Omega_2 \cdots \Omega_{K-1} \right\| \geq \frac{K-1}{K} \varepsilon \right\} \\
& \quad + \mathbb{P} \left\{ \prod_{k=1}^{K-1} M_k \left\| \hat{\Omega}_K^{(n)} - \Omega_K \right\| \geq \frac{1}{K} \varepsilon \right\} \\
& \leq \sum_{i=1}^{K-1} \xi_i \left(n, \frac{\varepsilon}{K M_K \prod_{k \neq i}^{K-1} M_k} \right) + \xi_K \left(n, \frac{\varepsilon}{K \prod_{k=1}^{K-1} M_k} \right) \\
& = \sum_{i=1}^K \xi_i \left(n, \frac{\varepsilon}{K \prod_{k \neq i}^K M_k} \right).
\end{aligned}$$

Thus, by induction, the proof is complete. \square

Remark A.1. In Lemma A.3, consider the special case where $\Omega_1 = \mathbf{0}$. Then the inequality (15) can be simplified as follows:

$$\left\| \hat{\Omega}_1^{(n)} \hat{\Omega}_2^{(n)} \cdots \hat{\Omega}_K^{(n)} - \mathbf{0} \right\| \leq \prod_{k=2}^K M_k \left\| \hat{\Omega}_1^{(n)} \right\|.$$

And we have the following simplified form:

$$\begin{aligned}
\mathbb{P} \left\{ \left\| \hat{\Omega}_1^{(n)} \hat{\Omega}_2^{(n)} \cdots \hat{\Omega}_K^{(n)} - \mathbf{0} \right\| \geq \varepsilon \right\} & \leq \mathbb{P} \left\{ \prod_{k=2}^K M_k \left\| \hat{\Omega}_1^{(n)} \right\| \geq \varepsilon \right\} \\
& \leq \xi_1 \left(n, \frac{\varepsilon}{\prod_{k=2}^K M_k} \right).
\end{aligned}$$

Lemma A.4 (Deviation Inequality for Minimum Eigenvalue of Projected Sample Covariance Matrix). Suppose Assumption 2 holds.

When $n \geq \max \left\{ \frac{qe^{\frac{3}{2}}}{K}, \frac{p^2(q+1)^2 K}{qK_0^2} \right\}$, the following inequality holds with probability at least $1 - \frac{3qe^{\frac{1}{2}}}{Kn}$:

$$\lambda_{\min} \left(\frac{1}{n} \mathbf{X}^\top \mathbf{P}_Z \mathbf{X} \right) \geq \lambda_z \left(\sigma_{\min}(\Theta) - \sqrt{\frac{2p(q+1)B_{\varepsilon_2}^2 \log(\frac{Kn}{q})}{\lambda_{\min}(\Sigma_z)n}} \right)^2 := \lambda_{\tilde{x}},$$

where $K := \frac{\lambda_{\min}(\Sigma_z)}{6B_{\varepsilon_2}^2}$, $K_0 := \frac{\lambda_{\min}(\Sigma_z)\sigma_{\min}^2(\Theta)}{2B_{\varepsilon_2}^2}$, $\Sigma_z := \mathbb{E}[zz^\top]$, $\mathbf{P}_Z := \mathbf{Z}(\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top$, and λ_z is a lower bound of $\lambda_{\min}(\frac{\mathbf{Z}^\top \mathbf{Z}}{n})$.

Proof. Let

$$\mathbf{E}_{\parallel} := \mathbf{P}_Z \mathbf{E}_2, \mathbf{E}_{\perp} := (\mathbf{I} - \mathbf{P}_Z) \mathbf{E}_2.$$

We have the following decomposition:

$$\mathbf{X} = \mathbf{Z}\Theta + \mathbf{E}_2 = \mathbf{Z}\Theta + \mathbf{E}_{\parallel} + \mathbf{E}_{\perp} = \mathbf{Z}(\Theta + (\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top \mathbf{E}_2) + \mathbf{E}_{\perp}.$$

Let $\Psi := (\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top \mathbf{E}_2$. Since $\mathbf{P}_Z \mathbf{E}_{\perp} = \mathbf{0}$, we have

$$\begin{aligned}
\mathbf{X}^\top \mathbf{P}_Z \mathbf{X} & = (\mathbf{Z}(\Theta + \Psi))^\top \mathbf{P}_Z (\mathbf{Z}(\Theta + \Psi)) \\
& = (\Theta + \Psi)^\top \mathbf{Z}^\top \mathbf{Z} (\Theta + \Psi).
\end{aligned}$$

We can now write:

$$\lambda_{\min} \left(\frac{1}{n} \mathbf{X}^\top \mathbf{P}_Z \mathbf{X} \right) = \lambda_{\min} \left((\boldsymbol{\Theta} + \boldsymbol{\Psi})^\top \frac{\mathbf{Z}^\top \mathbf{Z}}{n} (\boldsymbol{\Theta} + \boldsymbol{\Psi}) \right). \quad (16)$$

Note that in general, for a positive semi-definite matrix \mathbf{A} , we have

$$\lambda_{\min}(\mathbf{A}) = \min_{\mathbf{u} \neq \mathbf{0}} \frac{\mathbf{u}^\top \mathbf{A} \mathbf{u}}{\mathbf{u}^\top \mathbf{u}},$$

and

$$\begin{aligned} \lambda_{\min}(\mathbf{B}^\top \mathbf{A} \mathbf{B}) &= \min_{\mathbf{u} \neq \mathbf{0}} \frac{(\mathbf{B} \mathbf{u})^\top \mathbf{A} \mathbf{B} \mathbf{u}}{\mathbf{u}^\top \mathbf{u}} \\ &\geq \min_{\mathbf{u} \neq \mathbf{0}} \lambda_{\min}(\mathbf{A}) \frac{(\mathbf{B} \mathbf{u})^\top \mathbf{B} \mathbf{u}}{\mathbf{u}^\top \mathbf{u}} \\ &= \lambda_{\min}(\mathbf{A}) \lambda_{\min}(\mathbf{B}^\top \mathbf{B}). \end{aligned}$$

Thus, from Equation (16), with probability at least $1 - \xi$, we have:

$$\begin{aligned} \lambda_{\min} \left(\frac{1}{n} \mathbf{X}^\top \mathbf{P}_Z \mathbf{X} \right) &= \lambda_{\min} \left((\boldsymbol{\Theta} + \boldsymbol{\Psi})^\top \frac{\mathbf{Z}^\top \mathbf{Z}}{n} (\boldsymbol{\Theta} + \boldsymbol{\Psi}) \right) \\ &\geq \lambda_z \lambda_{\min} \left((\boldsymbol{\Theta} + \boldsymbol{\Psi})^\top (\boldsymbol{\Theta} + \boldsymbol{\Psi}) \right). \end{aligned}$$

It now remains to bound $\lambda_{\min} \left((\boldsymbol{\Theta} + \boldsymbol{\Psi})^\top (\boldsymbol{\Theta} + \boldsymbol{\Psi}) \right) = \sigma_{\min}^2(\boldsymbol{\Theta} + \boldsymbol{\Psi})$.

From (Hsu et al., 2014), for each $k \in [p]$ and any given $t > 1$, with sample size satisfying

$$n \geq \frac{6B_z^2(\log q + t)}{\lambda_{\min}(\boldsymbol{\Sigma}_z)}, \quad (17)$$

we have the following holds with probability at least $1 - 3e^{-t}$:

$$\|\boldsymbol{\Psi}_k\|_{\boldsymbol{\Sigma}_z}^2 = \|\hat{\boldsymbol{\Theta}}_k - \boldsymbol{\Theta}_k\|_{\boldsymbol{\Sigma}_z}^2 \leq \frac{B_{\epsilon_2}^2 (q + 2\sqrt{qt} + 2t)}{n} < \frac{B_{\epsilon_2}^2 [q + 2(q+1)t]}{n}.$$

Note that

$$\|\boldsymbol{\Psi}_k\|_{\boldsymbol{\Sigma}_z}^2 = \boldsymbol{\Psi}_k^\top \boldsymbol{\Sigma}_z \boldsymbol{\Psi}_k = \boldsymbol{\Psi}_k^\top \mathbf{U} \boldsymbol{\Lambda}_z \mathbf{U}^\top \boldsymbol{\Psi}_k = \sum_{i=1}^q \lambda_{z,i} (\mathbf{U}^\top \boldsymbol{\Psi}_k)_i^2,$$

and

$$\|\boldsymbol{\Psi}_k\| = \boldsymbol{\Psi}_k^\top \boldsymbol{\Psi}_k = \boldsymbol{\Psi}_k^\top \mathbf{U} \mathbf{U}^\top \boldsymbol{\Psi}_k = \sum_{i=1}^q (\mathbf{U}^\top \boldsymbol{\Psi}_k)_i^2.$$

We have

$$\lambda_{\min}(\boldsymbol{\Sigma}_z) \|\boldsymbol{\Psi}_k\|^2 \leq \|\boldsymbol{\Psi}_k\|_{\boldsymbol{\Sigma}_z}^2 \leq \lambda_{\max}(\boldsymbol{\Sigma}_z) \|\boldsymbol{\Psi}_k\|^2.$$

Then

$$\|\boldsymbol{\Psi}\| \leq \|\boldsymbol{\Psi}\|_F = \sqrt{\sum_{k=1}^p \|\boldsymbol{\Psi}_k\|^2} \leq \sqrt{\sum_{k=1}^p \frac{1}{\lambda_{\min}(\boldsymbol{\Sigma}_z)} \|\boldsymbol{\Psi}_k\|_{\boldsymbol{\Sigma}_z}^2} < \sqrt{\frac{pB_{\epsilon_2}^2 [q + 2(q+1)t]}{\lambda_{\min}(\boldsymbol{\Sigma}_z)n}}. \quad (18)$$

Hence, by Weyl's inequality, we have

$$\sigma_{\min}(\boldsymbol{\Theta} + \boldsymbol{\Psi}) \geq \sigma_{\min}(\boldsymbol{\Theta}) - \|\boldsymbol{\Psi}\| > \sigma_{\min}(\boldsymbol{\Theta}) - \sqrt{\frac{pB_{\epsilon_2}^2 [q + 2(q+1)t]}{\lambda_{\min}(\boldsymbol{\Sigma}_z)n}}, \quad (19)$$

where t is taken to be small enough such that the RHS ≥ 0 , i.e.

$$1 < t \leq \frac{K_0 n}{p(q+1)} - \frac{q}{2(q+1)}, \quad (20)$$

where $K_0 := \frac{\lambda_{\min}(\boldsymbol{\Sigma}_z) \sigma_{\min}^2(\boldsymbol{\Theta})}{2B_{\epsilon_2}^2}$. We now rewrite inequality Equation (19) in terms of n only. From condition Equation (17), for any given sample size n , the range for t is:

$$1 < t \leq Kn - \log q, \quad (21)$$

where $K := \frac{\lambda_{\min}(\boldsymbol{\Sigma}_z)}{6B_z^2}$. We take

$$t = \log(Kn) - \log q - \frac{1}{2} = \log\left(\frac{K}{q}n\right) - \frac{1}{2},$$

So that condition Equation (21) is satisfied when $n \geq \frac{qe^{\frac{3}{2}}}{K}$. To satisfy condition Equation (20), a sufficient condition is:

$$\log\left(\frac{K}{q}n\right) \leq \frac{K_0 n}{p(q+1)}.$$

Note that when $n \geq \frac{qe^{\frac{3}{2}}}{K}$, we also have:

$$\log\left(\frac{K}{q}n\right) \leq \sqrt{\frac{K}{q}n}.$$

So a sufficient condition to satisfy both Equation (20) and Equation (21) is:

$$n \geq \max\left\{\frac{qe^{\frac{3}{2}}}{K}, \frac{p^2(q+1)^2 K}{qK_0^2}\right\}.$$

Then the bound Equation (18) can be rewritten as:

$$\|\boldsymbol{\Psi}\| \leq \sqrt{\frac{pB_{\epsilon_2}^2 \left[q + 2(q+1) \left(\log\left(\frac{K}{q}n\right) - \frac{1}{2} \right) \right]}{\lambda_{\min}(\boldsymbol{\Sigma}_z)n}} < \sqrt{\frac{2p(q+1)B_{\epsilon_2}^2 \log\left(\frac{K}{q}n\right)}{\lambda_{\min}(\boldsymbol{\Sigma}_z)n}}. \quad (22)$$

Finally, from Equation (16),

$$\lambda_{\min}\left(\frac{1}{n}\mathbf{X}^\top \mathbf{P}_Z \mathbf{X}\right) \geq \lambda_z \left(\sigma_{\min}(\boldsymbol{\Theta}) - \sqrt{\frac{2p(q+1)B_{\epsilon_2}^2 \log\left(\frac{K}{q}n\right)}{\lambda_{\min}(\boldsymbol{\Sigma}_z)n}} \right)^2.$$

□

Proof of Theorem 2.1. We denote the observational values $(\mathbf{Z}, \mathbf{X}, \mathbf{Y}) = \{(z_i, \mathbf{x}_i, y_i)\}_{i=1}^n$, and $\boldsymbol{\epsilon}_1 = \{\epsilon_{1,i}\}_{i=1}^n$, $\boldsymbol{\epsilon}_2 = \{\epsilon_{2,i}\}_{i=1}^n$. The 2SLS estimator is given by:

$$\begin{aligned} \hat{\boldsymbol{\beta}}_{2\text{SLS}} &= \left(\hat{\boldsymbol{\Theta}}^\top \mathbf{Z}^\top \mathbf{Z} \hat{\boldsymbol{\Theta}} \right)^{-1} \hat{\boldsymbol{\Theta}}^\top \mathbf{Z}^\top \mathbf{Y} \\ &= \left[\left((\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top \mathbf{X} \right)^\top \mathbf{Z}^\top \mathbf{Z} (\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top \mathbf{X} \right]^{-1} \left((\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top \mathbf{X} \right)^\top \mathbf{Z}^\top \mathbf{Y} \\ &= (\mathbf{X}^\top \mathbf{Z} (\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Z} (\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top \mathbf{Y} \\ &= \boldsymbol{\beta} + (\mathbf{X}^\top \mathbf{Z} (\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Z} (\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top \boldsymbol{\epsilon}_1. \end{aligned} \quad (23)$$

Define constants $\lambda_z, \lambda_{\bar{x}} > 0$, such that the following event \mathcal{A} holds with probability at least $1 - \xi$:

$$\mathcal{A} = \left\{ \lambda_{\min}\left(\frac{\mathbf{Z}^\top \mathbf{Z}}{n}\right) \geq \lambda_z, \lambda_{\min}\left(\frac{\mathbf{X}^\top \mathbf{P}_Z \mathbf{X}}{n}\right) \geq \lambda_{\bar{x}} \right\}, \quad (24)$$

where $\lambda_{\min}(\cdot)$ denotes the smallest eigenvalue of a matrix, $\mathbf{P}_Z := \mathbf{Z}(\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top$ denotes the projection matrix. We will first assume the existence of such $\lambda_z, \lambda_{\bar{x}}$, with their values to be determined later.

We first consider the case when event \mathcal{A} is true. Let $\mathbf{Q}_{zz} := \mathbb{E}[\mathbf{z}\mathbf{z}^\top | \mathcal{A}]$, $\mathbf{Q}_{zx} := \mathbb{E}[\mathbf{z}\mathbf{x}^\top | \mathcal{A}]$, $\bar{\boldsymbol{\Omega}}_{zz} := \sum_{i=1}^n (\mathbf{z}_i \mathbf{z}_i^\top - \mathbf{Q}_{zz})$, $\bar{\boldsymbol{\Omega}}_{zx} := \sum_{i=1}^n (\mathbf{z}_i \mathbf{x}_i^\top - \mathbf{Q}_{zx})$, $\boldsymbol{\Omega}_{z\epsilon_1} := \sum_{i=1}^n \mathbf{z}_i \epsilon_{1,i}$.

Let $\bar{B}_{zz}, \bar{B}_{zx}, B_{zx}, B_{z\epsilon_1}$ be some upper bounds such that $\|\mathbf{z}_i \mathbf{z}_i^\top - \mathbf{Q}_{zz}\| \leq \bar{B}_{zz}$, $\|\mathbf{z}_i \mathbf{x}_i^\top - \mathbf{Q}_{zx}\| \leq \bar{B}_{zx}$, $\|\mathbf{z}_i \mathbf{x}_i^\top\| \leq B_{zx}$, $\|\mathbf{z}_i \epsilon_{1,i}\| \leq B_{z\epsilon_1}$ almost surely, for all $i = 1, \dots, n$. The existence of $\bar{B}_{zz}, \bar{B}_{zx}, B_{zx}, B_{z\epsilon_1}$ is guaranteed under Assumption 2(ii).

By Lemma A.1, we have:

$$\begin{aligned} \mathbb{P} \left\{ \left\| \frac{\mathbf{Z}^\top \mathbf{Z}}{n} - \mathbf{Q}_{zz} \right\| \geq \varepsilon \middle| \mathcal{A} \right\} &= \mathbb{P} \left\{ \left\| \frac{\sum_{i=1}^n \mathbf{z}_i \mathbf{z}_i^\top}{n} - \mathbf{Q}_{zz} \right\| \geq \varepsilon \middle| \mathcal{A} \right\} \\ &= \mathbb{P} \left\{ \left\| \sum_{i=1}^n (\mathbf{z}_i \mathbf{z}_i^\top - \mathbf{Q}_{zz}) \right\| \geq n\varepsilon \middle| \mathcal{A} \right\} \\ &\leq 2q \exp \left(-\frac{n^2 \varepsilon^2 / 2}{\nu(\bar{\boldsymbol{\Omega}}_{zz} | \mathcal{A}) + \bar{B}_{zz} n \varepsilon / 3} \right). \end{aligned} \quad (25)$$

Similarly,

$$\begin{aligned} \mathbb{P} \left\{ \left\| \frac{\mathbf{Z}^\top \mathbf{X}}{n} - \mathbf{Q}_{zx} \right\| \geq \varepsilon \middle| \mathcal{A} \right\} &= \mathbb{P} \left\{ \left\| \frac{\sum_{i=1}^n \mathbf{z}_i \mathbf{x}_i^\top}{n} - \mathbf{Q}_{zx} \right\| \geq \varepsilon \middle| \mathcal{A} \right\} \\ &= \mathbb{P} \left\{ \left\| \sum_{i=1}^n (\mathbf{z}_i \mathbf{x}_i^\top - \mathbf{Q}_{zx}) \right\| \geq n\varepsilon \middle| \mathcal{A} \right\} \\ &\leq (p+q) \exp \left(-\frac{n^2 \varepsilon^2 / 2}{\nu(\bar{\boldsymbol{\Omega}}_{zx} | \mathcal{A}) + \bar{B}_{zx} n \varepsilon / 3} \right). \end{aligned} \quad (26)$$

By Assumption 1(iii), the instrument \mathbf{z} is uncorrelated with the error term ϵ_1 , which implies $\mathbb{E}[\mathbf{z}\epsilon_1 | \mathcal{A}] = \mathbf{0}$. Applying Lemma A.1 again, we have:

$$\begin{aligned} \mathbb{P} \left\{ \left\| \frac{\mathbf{Z}^\top \boldsymbol{\epsilon}_1}{n} \right\| \geq \varepsilon \middle| \mathcal{A} \right\} &= \mathbb{P} \left\{ \left\| \frac{\sum_{i=1}^n \mathbf{z}_i \epsilon_{1,i}}{n} \right\| \geq \varepsilon \middle| \mathcal{A} \right\} \\ &= \mathbb{P} \left\{ \left\| \sum_{i=1}^n \mathbf{z}_i \epsilon_{1,i} \right\| \geq n\varepsilon \middle| \mathcal{A} \right\} \\ &\leq (q+1) \exp \left(-\frac{n^2 \varepsilon^2 / 2}{\nu(\boldsymbol{\Omega}_{z\epsilon_1} | \mathcal{A}) + B_{z\epsilon_1} n \varepsilon / 3} \right). \end{aligned} \quad (27)$$

With Lemma A.2 and (25), we have:

$$\begin{aligned} \mathbb{P} \left\{ \left\| n(\mathbf{Z}^\top \mathbf{Z})^{-1} - \mathbf{Q}_{ZZ}^{-1} \right\| \geq \varepsilon \middle| \mathcal{A} \right\} &\leq 2q \exp \left(-\frac{n^2 (\lambda_z^2 \varepsilon)^2 / 2}{\nu(\bar{\boldsymbol{\Omega}}_{zz} | \mathcal{A}) + \bar{B}_{zz} n (\lambda_z^2 \varepsilon) / 3} \right) \\ &= 2q \exp \left(-\frac{\lambda_z^4 n^2 \varepsilon^2 / 2}{\nu(\bar{\boldsymbol{\Omega}}_{zz} | \mathcal{A}) + \lambda_z^2 \bar{B}_{zz} n \varepsilon / 3} \right). \end{aligned} \quad (28)$$

Note that we have $\hat{\boldsymbol{\Theta}} = \boldsymbol{\Theta} + (\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top \boldsymbol{\epsilon}_2 := \boldsymbol{\Theta} + \boldsymbol{\Psi}$. Under event \mathcal{A} ,

$$\|\hat{\boldsymbol{\Theta}}\| = \|\boldsymbol{\Theta} + \boldsymbol{\Psi}\| \leq \|\boldsymbol{\Theta}\| + \|\boldsymbol{\Psi}\| \leq B_{\boldsymbol{\Theta}} + B_{\boldsymbol{\Psi}} := B_{\hat{\boldsymbol{\Theta}}}. \quad (29)$$

With Lemma A.3 (Remark A.1), combining (27)(29), we have:

$$\begin{aligned} \mathbb{P} \left\{ \left\| \frac{1}{n} \mathbf{X}^\top \mathbf{Z} (\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top \boldsymbol{\epsilon}_1 - \mathbf{0} \right\| \geq \varepsilon \middle| \mathcal{A} \right\} &\leq (q+1) \exp \left(-\frac{n^2 (\frac{\varepsilon}{B_{\hat{\boldsymbol{\Theta}}}})^2 / 2}{\nu(\boldsymbol{\Omega}_{z\epsilon_1} | \mathcal{A}) + B_{z\epsilon_1} n (\frac{\varepsilon}{B_{\hat{\boldsymbol{\Theta}}}}) / 3} \right) \\ &= (q+1) \exp \left(-\frac{n^2 \varepsilon^2 / 2}{B_{\hat{\boldsymbol{\Theta}}}^2 \nu(\boldsymbol{\Omega}_{z\epsilon_1} | \mathcal{A}) + B_{\hat{\boldsymbol{\Theta}}} B_{z\epsilon_1} n \varepsilon / 3} \right). \end{aligned} \quad (30)$$

Additionally, with Lemma A.3, combining (26)(28), we have:

$$\begin{aligned}
& \mathbb{P} \left\{ \left\| \frac{1}{n} \mathbf{X}^\top \mathbf{Z} (\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top \mathbf{X} - \mathbf{Q}_{zx}^\top \mathbf{Q}_{zz}^{-1} \mathbf{Q}_{zx} \right\| \geq \varepsilon \middle| \mathcal{A} \right\} \\
& \leq 2(p+q) \exp \left(-\frac{n^2 (\frac{\lambda_z \varepsilon}{3\bar{B}_{zx}})^2 / 2}{\nu(\bar{\Omega}_{zx} | \mathcal{A}) + \bar{B}_{zx} n (\frac{\lambda_z \varepsilon}{3\bar{B}_{zx}}) / 3} \right) + 2q \exp \left(-\frac{\lambda_z^4 n^2 (\frac{\varepsilon}{3\bar{B}_{zx}^2})^2 / 2}{\nu(\bar{\Omega}_{zz} | \mathcal{A}) + \lambda_z^2 \bar{B}_{zz} n (\frac{\varepsilon}{3\bar{B}_{zx}^2}) / 3} \right) \\
& = 2(p+q) \exp \left(-\frac{\lambda_z^2 n^2 \varepsilon^2 / 2}{9\bar{B}_{zx}^2 \nu(\bar{\Omega}_{zx} | \mathcal{A}) + \lambda_z \bar{B}_{zx} \bar{B}_{zx} n \varepsilon} \right) + 2q \exp \left(-\frac{\lambda_z^4 n^2 \varepsilon^2 / 2}{9\bar{B}_{zx}^4 \nu(\bar{\Omega}_{zz} | \mathcal{A}) + \lambda_z^2 \bar{B}_{zx}^2 \bar{B}_{zz} n \varepsilon} \right).
\end{aligned}$$

Applying Lemma A.2 again, we have:

$$\begin{aligned}
& \mathbb{P} \left\{ \left\| n (\mathbf{X}^\top \mathbf{Z} (\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top \mathbf{X})^{-1} - (\mathbf{Q}_{ZX}^\top \mathbf{Q}_{zz}^{-1} \mathbf{Q}_{zx})^{-1} \right\| \geq \varepsilon \middle| \mathcal{A} \right\} \\
& \leq 2(p+q) \exp \left(-\frac{\lambda_z^2 n^2 (\lambda_{\bar{x}}^2 \varepsilon)^2 / 2}{9\bar{B}_{zx}^2 \nu(\bar{\Omega}_{zx} | \mathcal{A}) + \lambda_z \bar{B}_{zx} \bar{B}_{zx} n (\lambda_{\bar{x}}^2 \varepsilon)} \right) + 2q \exp \left(-\frac{\lambda_z^4 n^2 (\lambda_{\bar{x}}^2 \varepsilon)^2 / 2}{9\bar{B}_{zx}^4 \nu(\bar{\Omega}_{zz} | \mathcal{A}) + \lambda_z^2 \bar{B}_{zx}^2 \bar{B}_{zz} n (\lambda_{\bar{x}}^2 \varepsilon)} \right) \\
& = 2(p+q) \exp \left(-\frac{\lambda_z^2 \lambda_{\bar{x}}^4 n^2 \varepsilon^2 / 2}{9\bar{B}_{zx}^2 \nu(\bar{\Omega}_{zx} | \mathcal{A}) + \lambda_z \lambda_{\bar{x}}^2 \bar{B}_{zx} \bar{B}_{zx} n \varepsilon} \right) + 2q \exp \left(-\frac{\lambda_z^4 \lambda_{\bar{x}}^4 n^2 \varepsilon^2 / 2}{9\bar{B}_{zx}^4 \nu(\bar{\Omega}_{zz} | \mathcal{A}) + \lambda_z^2 \lambda_{\bar{x}}^2 \bar{B}_{zx}^2 \bar{B}_{zz} n \varepsilon} \right). \tag{31}
\end{aligned}$$

Therefore, we have shown that under event \mathcal{A} ,

$$n (\mathbf{X}^\top \mathbf{Z} (\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top \mathbf{X})^{-1} \xrightarrow{\mathbb{P}} (\mathbf{Q}_{ZX}^\top \mathbf{Q}_{zz}^{-1} \mathbf{Q}_{zx})^{-1}.$$

From equation (23), combining (30) and (31) with Lemma A.3 (Remark A.1), we have:

$$\begin{aligned}
& \mathbb{P} \left\{ \left\| \hat{\beta}_{2\text{SLS}} - \beta \right\| \geq \varepsilon \middle| \mathcal{A} \right\} \\
& = \mathbb{P} \left\{ \left\| (\mathbf{X}^\top \mathbf{Z} (\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Z} (\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top \boldsymbol{\varepsilon}_1 - \mathbf{0} \right\| \geq \varepsilon \middle| \mathcal{A} \right\} \\
& \leq (q+1) \exp \left(-\frac{n^2 (\lambda_{\bar{x}} \varepsilon)^2 / 2}{\bar{B}_{\hat{\Theta}}^2 \nu(\bar{\Omega}_{z\epsilon_1} | \mathcal{A}) + \bar{B}_{\hat{\Theta}} \bar{B}_{z\epsilon_1} n (\lambda_{\bar{x}} \varepsilon) / 3} \right) \\
& = (q+1) \exp \left(-\frac{\lambda_{\bar{x}}^2 n^2 \varepsilon^2 / 2}{\bar{B}_{\hat{\Theta}}^2 \nu(\bar{\Omega}_{z\epsilon_1} | \mathcal{A}) + \lambda_{\bar{x}} \bar{B}_{\hat{\Theta}} \bar{B}_{z\epsilon_1} n \varepsilon / 3} \right).
\end{aligned}$$

For the second part of the theorem, let $c := \left(\frac{3\bar{B}_{\hat{\Theta}} \nu(\bar{\Omega}_{z\epsilon_1} | \mathcal{A})}{\lambda_{\bar{x}} \bar{B}_{z\epsilon_1} n} \right)^2$, we have:

$$\begin{aligned}
& \mathbb{E} \left[\left\| \text{clip}_{B_\beta}(\hat{\beta}_{2\text{SLS}}) - \beta \right\|^2 \right] \\
& = \mathbb{E} \left[\left\| \text{clip}_{B_\beta}(\hat{\beta}_{2\text{SLS}}) - \beta \right\|^2 \middle| \mathcal{A} \right] \mathbb{P} \{ \mathcal{A} \} + \mathbb{E} \left[\left\| \text{clip}_{B_\beta}(\hat{\beta}_{2\text{SLS}}) - \beta \right\|^2 \middle| \mathcal{A}^c \right] \mathbb{P} \{ \mathcal{A}^c \} \\
& \leq \mathbb{E} \left[\left\| \hat{\beta}_{2\text{SLS}} - \beta \right\|^2 \middle| \mathcal{A} \right] \mathbb{P} \{ \mathcal{A} \} + \mathbb{E} \left[\left\| \text{clip}_{B_\beta}(\hat{\beta}_{2\text{SLS}}) - \beta \right\|^2 \middle| \mathcal{A}^c \right] \mathbb{P} \{ \mathcal{A}^c \} \\
& \leq \mathbb{E} \left[\left\| \hat{\beta}_{2\text{SLS}} - \beta \right\|^2 \middle| \mathcal{A} \right] + \mathbb{E} \left[\left\| \text{clip}_{B_\beta}(\hat{\beta}_{2\text{SLS}}) - \beta \right\|^2 \middle| \mathcal{A}^c \right] \cdot \xi,
\end{aligned} \tag{32}$$

where

$$\mathbb{E} \left[\left\| \text{clip}_{B_\beta}(\hat{\beta}_{2\text{SLS}}) - \beta \right\|^2 \middle| \mathcal{A}^c \right] \leq 4B_\beta^2, \tag{33}$$

and

$$\begin{aligned}
& \mathbb{E} \left[\|\hat{\beta}_{2\text{SLS}} - \beta\|^2 \middle| \mathcal{A} \right] \\
&= \int_0^\infty \mathbb{P} \left\{ \|\hat{\beta}_{2\text{SLS}} - \beta\|^2 \geq \varepsilon \middle| \mathcal{A} \right\} d\varepsilon \\
&= \int_0^\infty \mathbb{P} \left\{ \|\hat{\beta}_{2\text{SLS}} - \beta\| \geq \sqrt{\varepsilon} \middle| \mathcal{A} \right\} d\varepsilon \\
&\leq \int_0^\infty (q+1) \exp \left(-\frac{\lambda_{\bar{x}}^2 n^2 \varepsilon / 2}{B_{\hat{\Theta}}^2 \nu(\Omega_{z_{\epsilon_1}} | \mathcal{A}) + \lambda_{\bar{x}} B_{\hat{\Theta}} B_{z_{\epsilon_1}} n \sqrt{\varepsilon} / 3} \right) d\varepsilon \\
&\leq (q+1) \left[\int_0^c \exp \left(-\frac{\lambda_{\bar{x}}^2 n^2 \varepsilon / 2}{2B_{\hat{\Theta}}^2 \nu(\Omega_{z_{\epsilon_1}} | \mathcal{A})} \right) d\varepsilon + \int_c^\infty \exp \left(-\frac{\lambda_{\bar{x}} n \sqrt{\varepsilon} / 2}{2B_{\hat{\Theta}} B_{z_{\epsilon_1}} / 3} \right) d\varepsilon \right] \\
&= (q+1) \left[\frac{4B_{\hat{\Theta}}^2 \nu(\Omega_{z_{\epsilon_1}} | \mathcal{A})}{\lambda_{\bar{x}}^2 n^2} \left(1 - \exp \left(-\frac{9\nu(\Omega_{z_{\epsilon_1}} | \mathcal{A})}{4B_{z_{\epsilon_1}}^2} \right) \right) + \left(\frac{8B_{\hat{\Theta}}^2 \nu(\Omega_{z_{\epsilon_1}} | \mathcal{A})}{\lambda_{\bar{x}}^2 n^2} + \frac{32B_{\hat{\Theta}}^2 B_{z_{\epsilon_1}}^2}{9\lambda_{\bar{x}}^2 n^2} \right) \exp \left(-\frac{9\nu(\Omega_{z_{\epsilon_1}} | \mathcal{A})}{4B_{z_{\epsilon_1}}^2} \right) \right] \\
&\leq (q+1) \left[\frac{4B_{\hat{\Theta}}^2 \nu(\Omega_{z_{\epsilon_1}} | \mathcal{A})}{\lambda_{\bar{x}}^2 n^2} + \frac{8B_{\hat{\Theta}}^2 \nu(\Omega_{z_{\epsilon_1}} | \mathcal{A})}{\lambda_{\bar{x}}^2 n^2} + \frac{32B_{\hat{\Theta}}^2 B_{z_{\epsilon_1}}^2}{9\lambda_{\bar{x}}^2 n^2} \right] \\
&= \frac{(q+1)B_{\hat{\Theta}}^2}{\lambda_{\bar{x}}^2 n^2} \left[12\nu(\Omega_{z_{\epsilon_1}} | \mathcal{A}) + \frac{32B_z^2 B_{\epsilon_1}^2}{9} \right]. \tag{34}
\end{aligned}$$

Note that we further have the following bound:

$$\begin{aligned}
\nu(\Omega_{z_{\epsilon_1}} | \mathcal{A}) &= \max \left\{ \left\| \mathbb{E} \left[\left(\sum_{i=1}^n z_{i\epsilon_{1,i}} \right)^\top \left(\sum_{j=1}^n z_{j\epsilon_{1,j}} \right) \middle| \mathcal{A} \right] \right\|, \left\| \mathbb{E} \left[\left(\sum_{i=1}^n z_{i\epsilon_{1,i}} \right) \left(\sum_{j=1}^n z_{j\epsilon_{1,j}} \right)^\top \middle| \mathcal{A} \right] \right\| \right\} \\
&= \max \left\{ \left\| \mathbb{E} \left[\sum_{i=1}^n \epsilon_{1,i}^2 z_i^\top z_i \middle| \mathcal{A} \right] \right\|, \left\| \mathbb{E} \left[\sum_{i=1}^n \epsilon_{1,i} z_i z_i^\top \middle| \mathcal{A} \right] \right\| \right\} \\
&\leq nB_z^2 \sigma_1^2. \tag{35}
\end{aligned}$$

It now remains to determine λ_z , $\lambda_{\bar{x}}$, and $B_{\hat{\Theta}}$.

From Theorem 4.6.1 of (Vershynin, 2018), when $n \geq c^2 B_z^4 (\sqrt{q} + \sqrt{t})^2$, with probability at least $1 - 2e^{-t}$, we have:

$$\lambda_{\min} \left(\frac{1}{n} \mathbf{Z}^\top \mathbf{Z} \right) \geq \lambda_{\min}(\Sigma_z) \left(1 - \frac{cB_z^2 (\sqrt{q} + \sqrt{t})}{\sqrt{n}} \right)^2,$$

where c is an absolute constant. We rewrite the theorem by taking $t = \log\left(\frac{K}{q}n\right) - \frac{1}{2}$ (similar to the proof in Lemma A.4). Then we need the following condition to be satisfied:

$$n \geq c^2 B_z^4 \left(\sqrt{q} + \sqrt{\log\left(\frac{K}{q}n\right) - \frac{1}{2}} \right)^2. \tag{36}$$

We can bound the RHS of Equation (36) as follows:

$$\begin{aligned}
c^2 B_z^4 \left(\sqrt{q} + \sqrt{\log\left(\frac{K}{q}n\right) - \frac{1}{2}} \right)^2 &\leq 2c^2 B_z^4 \left(q + \log\left(\frac{K}{q}n\right) - \frac{1}{2} \right) \\
&\leq \frac{n}{2} + 2c^2 B_z^4 \left(q + \log\left(\frac{4c^2 B_z^4 K}{q}\right) - \frac{3}{2} \right),
\end{aligned}$$

where the second line follows from the inequality $\log(x) \leq \frac{x}{C} + \log(C) - 1$, with $x = \frac{K}{q}n$ and $C = \frac{4c^2 B_z^4 K}{q}$. So a sufficient condition for Equation (36) to hold is:

$$n \geq 4c^2 B_z^4 \left(q + \log \left(\frac{4c^2 B_z^4 K}{q} \right) - \frac{3}{2} \right). \quad (37)$$

With condition Equation (37), we have the following bound holds with probability at least $1 - \frac{2qe^{\frac{1}{2}}}{Kn}$:

$$\lambda_{\min} \left(\frac{1}{n} \mathbf{Z}^\top \mathbf{Z} \right) \geq \lambda_{\min}(\boldsymbol{\Sigma}_z) \left(1 - \frac{cB_z^2 \left(\sqrt{q} + \sqrt{\log \left(\frac{K}{q}n \right) - \frac{1}{2}} \right)}{\sqrt{n}} \right)^2 := \lambda_z. \quad (38)$$

Furthermore, with Lemma A.4, when $n \geq \max \left\{ \frac{qe^{\frac{3}{2}}}{K}, \frac{p^2(q+1)^2 K}{qK_0^2} \right\}$, we have the following holds with probability at least $1 - \frac{5qe^{\frac{1}{2}}}{Kn}$:

$$\begin{aligned} & \lambda_{\min} \left(\frac{1}{n} \mathbf{X}^\top \mathbf{P}_Z \mathbf{X} \right) \\ & \geq \lambda_z \left(\sigma_{\min}(\boldsymbol{\Theta}) - \sqrt{\frac{2p(q+1)B_{\epsilon_2}^2 \log \left(\frac{K}{q}n \right)}{\lambda_{\min}(\boldsymbol{\Sigma}_z)n}} \right)^2 \\ & = \lambda_{\min}(\boldsymbol{\Sigma}_z) \left(1 - \frac{cB_z^2 \left(\sqrt{q} + \sqrt{\log \left(\frac{K}{q}n \right) - \frac{1}{2}} \right)}{\sqrt{n}} \right)^2 \left(\sigma_{\min}(\boldsymbol{\Theta}) - \sqrt{\frac{2p(q+1)B_{\epsilon_2}^2 \log \left(\frac{K}{q}n \right)}{\lambda_{\min}(\boldsymbol{\Sigma}_z)n}} \right)^2 \\ & := \lambda_{\bar{x}}. \end{aligned} \quad (39)$$

From Equation (22) and Equation (29), we have:

$$B_{\hat{\Theta}} = B_{\Theta} + \sqrt{\frac{2p(q+1)B_{\epsilon_2}^2 \log \left(\frac{K}{q}n \right)}{\lambda_{\min}(\boldsymbol{\Sigma}_z)n}}. \quad (40)$$

With $\xi = \frac{5qe^{\frac{1}{2}}}{Kn}$, putting together Equations (35)(38)(39)(40) into Equation (34), and Equations (33)(34) into Equation (32), we have:

$$\begin{aligned} \mathbb{E} \left[\|\text{clip}_{B_\beta}(\hat{\beta}_{2\text{SLS}}) - \beta\|^2 \right] & \leq \mathbb{E} \left[\|\hat{\beta}_{2\text{SLS}} - \beta\|^2 \mid \mathcal{A} \right] + 4B_\beta^2 \xi \\ & \leq \frac{(q+1)B_{\hat{\Theta}}^2}{\lambda_x^2 n^2} \left[12nB_z^2 \sigma_1^2 + \frac{32B_z^2 B_{\epsilon_1}^2}{9} \right] + \frac{20qe^{\frac{1}{2}} B_\beta^2}{Kn} \\ & \leq \mathcal{O} \left(\frac{q}{n} \left(\frac{B_\beta^2}{K} + C^2(n)\sigma_1^2 \right) \right), \end{aligned} \quad (41)$$

where

$$C(n) := \frac{\left(B_{\Theta} + \sqrt{\frac{2p(q+1)B_{\epsilon_2}^2 \log \left(\frac{K}{q}n \right)}{\lambda_{\min}(\boldsymbol{\Sigma}_z)n}} \right) B_z}{\lambda_{\min}(\boldsymbol{\Sigma}_z) \left(1 - \frac{cB_z^2 \left(\sqrt{q} + \sqrt{\log \left(\frac{K}{q}n \right) - \frac{1}{2}} \right)}{\sqrt{n}} \right)^2 \left(\sigma_{\min}(\boldsymbol{\Theta}) - \sqrt{\frac{2p(q+1)B_{\epsilon_2}^2 \log \left(\frac{K}{q}n \right)}{\lambda_{\min}(\boldsymbol{\Sigma}_z)n}} \right)^2}, \quad (42)$$

thus completing the proof. \square

B PROOFS FOR SECTION 3

B.1 PROOF OF THEOREM 3.1

Lemma B.1. Suppose $\{\Omega^{(1)}, \dots, \Omega^{(t)}, \dots\}$ is a $d \times d$ -matrix sequence decaying with exponential rate r , i.e. for some constant $c > 0$ and $0 < r < 1$,

$$\left\| \Omega^{(t)} \right\|_F \leq cr^t.$$

Then for any $\varepsilon > 0$, there exists a finite constant:

$$T_0 = \left\lceil \log_r \frac{(1-r)(\varepsilon/d)}{c(1+(1-r)(\varepsilon/d))} \right\rceil,$$

such that

$$\left\| \prod_{t=T_0}^{\infty} (I + \Omega^{(t)}) - I \right\|_F < \varepsilon,$$

and hence

$$\left\| \prod_{t=T_0}^{\infty} (I + \Omega^{(t)}) \right\|_F < \sqrt{d} + \varepsilon.$$

Proof. By definition,

$$\left\| \Omega^{(k)} \right\|_F = \sqrt{\sum_{i,j=1}^p \Omega_{ij}^{(k)2}} \leq cr^k,$$

which implies:

$$\left| \Omega_{ij}^{(k)} \right| \leq cr^k, \quad \forall i, j, k.$$

Consider the product of any two matrices. By sub-multiplicativity,

$$\left\| \Omega^{(k)} \Omega^{(l)} \right\|_F \leq \left\| \Omega^{(k)} \right\|_F \left\| \Omega^{(l)} \right\|_F \leq c^2 r^{k+l},$$

which implies:

$$\left| \left[\Omega^{(k)} \Omega^{(l)} \right]_{ij} \right| \leq c^2 r^{k+l}, \quad \forall i, j, k, l.$$

Similarly, for the product of any number of matrices:

$$\left| \left[\Omega^{(k_1)} \Omega^{(k_2)} \dots \Omega^{(k_n)} \right]_{ij} \right| \leq c^n r^{k_1+k_2+\dots+k_n}, \quad \forall i, j, k_1, \dots, k_n.$$

Thus

$$\begin{aligned} & \left\| \prod_{t=t_1}^{t_2} (I + \Omega^{(t)}) - I \right\|_F \\ &= \left\| (I + \Omega^{(t_1)}) (I + \Omega^{(t_1+1)}) \dots (I + \Omega^{(t_2)}) - I \right\|_F \\ &= \left\| \sum_{t_1 \leq k \leq t_2} \Omega^{(k)} + \sum_{t_1 \leq k < l \leq t_2} \Omega^{(k)} \Omega^{(l)} + \dots + \Omega^{(t_1)} \Omega^{(t_1+1)} \dots \Omega^{(t_2)} \right\|_F \\ &\leq \left\| \sum_{t_1 \leq k \leq t_2} cr^k \mathbf{1} \mathbf{1}^\top + \sum_{t_1 \leq k < l \leq t_2} c^2 r^{k+l} \mathbf{1} \mathbf{1}^\top + \dots + c^{t_2-t_1+1} r^{t_1+\dots+t_2} \mathbf{1} \mathbf{1}^\top \right\|_F. \end{aligned} \tag{43}$$

Note that the last inequality can be checked by comparing matrix elements of both sides. For any $\varepsilon > 0$, we take $T_0 = \lceil \log_r \frac{(1-r)(\varepsilon/d)}{c(1+(1-r)(\varepsilon/d))} \rceil$. Consider $t_1 = T_0$ and $t_2 \rightarrow \infty$ in (43). For notation convenience, let

$$\Xi := \sum_{T_0 \leq k} cr^k \mathbf{1}\mathbf{1}^\top + \sum_{T_0 \leq k < l} c^2 r^{k+l} \mathbf{1}\mathbf{1}^\top + \sum_{T_0 \leq k < l < m} c^3 r^{k+l+m} \mathbf{1}\mathbf{1}^\top + \dots$$

Then

$$\begin{aligned} \Xi_{ij} &= \sum_{T_0 \leq k} cr^k + \sum_{T_0 \leq k < l} c^2 r^{k+l} + \sum_{T_0 \leq k < l < m} c^3 r^{k+l+m} + \dots \\ &< c \sum_{k \geq T_0} r^k + c^2 r^{2T_0} \sum_{k \geq T_0} r^k + c^3 r^{3T_0} \sum_{k \geq T_0} r^k + \dots \\ &= \frac{cr^{T_0}}{1-r} + \frac{c^2 r^{2T_0}}{1-r} + \frac{c^3 r^{3T_0}}{1-r} + \dots \\ &= \frac{cr^{T_0}}{(1-r)(1-cr^{T_0})} \\ &\leq \frac{\varepsilon}{d}. \end{aligned}$$

Thus

$$\left\| \prod_{t=T_0}^{\infty} (\mathbf{I} + \Omega^{(t)}) - \mathbf{I} \right\|_F = \|\Xi\|_F = \sqrt{\sum_{i,j=1}^d \Xi_{ij}^2} \leq \varepsilon.$$

Hence completes the proof. \square

Proof of Theorem 3.1. In the following proof, we treat $\mathbf{Z}, \mathbf{X}, \mathbf{Y}$ as deterministic matrices.

We begin by checking the inner loop (7a):

$$\begin{aligned} \Theta^{(t)} - \hat{\Theta} &= \Theta^{(t-1)} - \hat{\Theta} - \eta \mathbf{Z}^\top (\mathbf{Z} \Theta^{(t-1)} - \mathbf{X}) \\ &= (\mathbf{I} - \eta \mathbf{Z}^\top \mathbf{Z}) (\Theta^{(t-1)} - \hat{\Theta}) + \eta \mathbf{Z}^\top (\mathbf{X} - \mathbf{Z} \hat{\Theta}) \\ &= (\mathbf{I} - \eta \mathbf{Z}^\top \mathbf{Z})^2 (\Theta^{(t-2)} - \hat{\Theta}) + \eta \mathbf{Z}^\top (\mathbf{X} - \mathbf{Z} \hat{\Theta}) + \eta (\mathbf{I} - \eta \mathbf{Z}^\top \mathbf{Z}) \mathbf{Z}^\top (\mathbf{X} - \mathbf{Z} \hat{\Theta}) \\ &\quad \vdots \\ &= (\mathbf{I} - \eta \mathbf{Z}^\top \mathbf{Z})^t (\Theta^{(0)} - \hat{\Theta}) + \sum_{i=0}^{t-1} \eta (\mathbf{I} - \eta \mathbf{Z}^\top \mathbf{Z})^{t-1-i} \mathbf{Z}^\top (\mathbf{X} - \mathbf{Z} \hat{\Theta}) \\ &= (\mathbf{I} - \eta \mathbf{Z}^\top \mathbf{Z})^t (\Theta^{(0)} - \hat{\Theta}) + \eta [\mathbf{I} - (\mathbf{I} - \eta \mathbf{Z}^\top \mathbf{Z})^t] (\eta \mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top (\mathbf{X} - \mathbf{Z} \hat{\Theta}) \\ &= (\mathbf{I} - \eta \mathbf{Z}^\top \mathbf{Z})^t (\Theta^{(0)} - \hat{\Theta}) + [\mathbf{I} - (\mathbf{I} - \eta \mathbf{Z}^\top \mathbf{Z})^t] [(\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top \mathbf{X} - \hat{\Theta}] \\ &= (\mathbf{I} - \eta \mathbf{Z}^\top \mathbf{Z})^t (\Theta^{(0)} - \hat{\Theta}). \end{aligned}$$

With learning rate $0 < \eta < \frac{2}{\sigma_{\max}^2(\mathbf{Z})}$, let $\kappa(\eta) := \rho(\mathbf{I} - \eta \mathbf{Z}^\top \mathbf{Z})$, where $\rho(\cdot)$ denotes the spectral radius. Then it follows that $0 < \kappa(\eta) < 1$. We have:

$$\begin{aligned} \|\Theta^{(t)} - \hat{\Theta}\| &= \|(\mathbf{I} - \eta \mathbf{Z}^\top \mathbf{Z})^t (\Theta^{(0)} - \hat{\Theta})\| \\ &\leq \|(\mathbf{I} - \eta \mathbf{Z}^\top \mathbf{Z})^t\| \|\Theta^{(0)} - \hat{\Theta}\| \\ &\leq \kappa(\eta)^t \|\Theta^{(0)} - \hat{\Theta}\| \\ &= \mathcal{O}(\kappa(\eta)^t). \end{aligned}$$

Thus $\{\Theta^{(t)}\}$ converges to $\hat{\Theta}$ exponentially with rate $\kappa(\eta)$.

For the outer loop (7b), we have:

$$\begin{aligned}
\beta^{(t)} - \hat{\beta}_{2\text{SLS}} &= \beta^{(t-1)} - \hat{\beta}_{2\text{SLS}} - \alpha \Theta^{(t-1)\top} \mathbf{Z}^\top \left(\mathbf{Z} \Theta^{(t-1)} \beta^{(t-1)} - \mathbf{Y} \right) \\
&= \left(\mathbf{I} - \alpha \Theta^{(t-1)\top} \mathbf{Z}^\top \mathbf{Z} \Theta^{(t-1)} \right) \left(\beta^{(t-1)} - \hat{\beta}_{2\text{SLS}} \right) + \alpha \Theta^{(t-1)\top} \mathbf{Z}^\top \left(\mathbf{Y} - \mathbf{Z} \Theta^{(t-1)} \hat{\beta}_{2\text{SLS}} \right) \\
&\quad \vdots \\
&= \underbrace{\prod_{i=0}^{t-1} \left(\mathbf{I} - \alpha \Theta^{(i)\top} \mathbf{Z}^\top \mathbf{Z} \Theta^{(i)} \right)}_{\Delta_1 \beta^{(t)}} \left(\beta^{(0)} - \hat{\beta}_{2\text{SLS}} \right) \\
&\quad + \underbrace{\sum_{i=0}^{t-1} \alpha \left[\prod_{j=i+1}^{t-1} \left(\mathbf{I} - \alpha \Theta^{(j)\top} \mathbf{Z}^\top \mathbf{Z} \Theta^{(j)} \right) \right]}_{\Delta_2 \beta^{(t)}} \Theta^{(i)\top} \mathbf{Z}^\top \left(\mathbf{Y} - \mathbf{Z} \Theta^{(i)} \hat{\beta}_{2\text{SLS}} \right).
\end{aligned} \tag{44}$$

To simplify notations, let

$$\begin{aligned}
\mathbf{R}^{(t)} &:= \Theta^{(t)} - \hat{\Theta} = \left(\mathbf{I} - \eta \mathbf{Z}^\top \mathbf{Z} \right)^t \left(\Theta^{(0)} - \hat{\Theta} \right), \\
\mathbf{V}^{(t)} &:= \left(\mathbf{I} - \alpha \hat{\Theta}^\top \mathbf{Z}^\top \mathbf{Z} \hat{\Theta} \right)^t, \\
\mathbf{W}^{(t)} &:= \mathbf{R}^{(t)\top} \mathbf{Z}^\top \mathbf{Z} \hat{\Theta} + \hat{\Theta}^\top \mathbf{Z}^\top \mathbf{Z} \mathbf{R}^{(t)} + \mathbf{R}^{(t)\top} \mathbf{Z}^\top \mathbf{Z} \mathbf{R}^{(t)}.
\end{aligned}$$

With learning rates $0 < \alpha < \frac{2}{\sigma_{\max}^2(\mathbf{Z}\hat{\Theta})}$, $0 < \eta < \frac{2}{\sigma_{\max}^2(\mathbf{Z})}$, let $\gamma(\alpha) := \rho \left(\mathbf{I} - \alpha \hat{\Theta}^\top \mathbf{Z}^\top \mathbf{Z} \hat{\Theta} \right)$. Then it follows that $0 < \gamma(\alpha) < 1$. We have:

$$\left\| \mathbf{R}^{(t)} \right\| \leq \kappa(\eta)^t \left\| \Theta^{(0)} - \hat{\Theta} \right\|, \tag{45}$$

$$\left\| \mathbf{V}^{(t)} \right\| \leq \gamma(\alpha)^t,$$

and

$$\begin{aligned}
\left\| \mathbf{W}^{(t)} \right\| &= \left\| \mathbf{R}^{(t)\top} \mathbf{Z}^\top \mathbf{Z} \hat{\Theta} + \hat{\Theta}^\top \mathbf{Z}^\top \mathbf{Z} \mathbf{R}^{(t)} + \mathbf{R}^{(t)\top} \mathbf{Z}^\top \mathbf{Z} \mathbf{R}^{(t)} \right\| \\
&\leq 2 \left\| \hat{\Theta}^\top \mathbf{Z}^\top \mathbf{Z} \right\| \left\| \mathbf{R}^{(t)} \right\| + \left\| \mathbf{Z}^\top \mathbf{Z} \right\| \left\| \mathbf{R}^{(t)} \right\|^2 \\
&\leq 2\kappa(\eta)^t \left\| \hat{\Theta}^\top \mathbf{Z}^\top \mathbf{Z} \right\| \left\| \Theta^{(0)} - \hat{\Theta} \right\| + \kappa(\eta)^{2t} \left\| \mathbf{Z}^\top \mathbf{Z} \right\| \left\| \Theta^{(0)} - \hat{\Theta} \right\|^2 \\
&\leq \kappa(\eta)^t \left(2 \left\| \hat{\Theta}^\top \mathbf{Z}^\top \mathbf{Z} \right\| \left\| \Theta^{(0)} - \hat{\Theta} \right\| + \left\| \mathbf{Z}^\top \mathbf{Z} \right\| \left\| \Theta^{(0)} - \hat{\Theta} \right\|^2 \right) \\
&= \mathcal{O}(\kappa(\eta)^t).
\end{aligned}$$

Then from Equation (44), we have:

$$\begin{aligned}
\Delta_1 \beta^{(t)} &= \prod_{i=0}^{t-1} \left(I - \alpha \Theta^{(i)\top} Z^\top Z \Theta^{(i)} \right) \left(\beta^{(0)} - \hat{\beta}_{2\text{SLS}} \right) \\
&= \prod_{i=0}^{t-1} \left[I - \alpha \hat{\Theta}^\top Z^\top Z \hat{\Theta} - \alpha \left(R^{(i)\top} Z^\top Z \hat{\Theta} + \hat{\Theta}^\top Z^\top Z R^{(i)} + R^{(i)\top} Z^\top Z R^{(i)} \right) \right] \left(\beta^{(0)} - \hat{\beta}_{2\text{SLS}} \right) \\
&= \prod_{i=0}^{t-1} \left[I - \alpha \hat{\Theta}^\top Z^\top Z \hat{\Theta} - \alpha W^{(i)} \right] \left(\beta^{(0)} - \hat{\beta}_{2\text{SLS}} \right) \\
&= \left(I - \alpha \hat{\Theta}^\top Z^\top Z \hat{\Theta} \right)^t \prod_{i=0}^{t-1} \left[I - \alpha \left(I - \alpha \hat{\Theta}^\top Z^\top Z \hat{\Theta} \right)^{-1} W^{(i)} \right] \left(\beta^{(0)} - \hat{\beta}_{2\text{SLS}} \right) \\
&= V^{(t)} \prod_{i=0}^{t-1} \left[I - \alpha \left(I - \alpha \hat{\Theta}^\top Z^\top Z \hat{\Theta} \right)^{-1} W^{(i)} \right] \left(\beta^{(0)} - \hat{\beta}_{2\text{SLS}} \right).
\end{aligned}$$

We denote $\Psi := \alpha \left(I - \alpha \hat{\Theta}^\top Z^\top Z \hat{\Theta} \right)^{-1}$. By Lemma B.1, we take $\varepsilon = 1$, c_0 be a constant such that $\|W^{(t)}\|_F \leq c_0 \kappa(\eta)^t$, and $T_0 = \lceil \log_{\kappa(\eta)} \frac{(1-\kappa(\eta))}{\|\Psi\|_F c_0 (p+(1-\kappa(\eta)))} \rceil$.

Then we have:

$$\left\| \prod_{i=T_0}^{t-1} \left(I - \Psi W^{(i)} \right) \right\| \leq \left\| \prod_{i=T_0}^{t-1} \left(I - \Psi W^{(i)} \right) \right\|_F < \sqrt{p} + 1. \quad (46)$$

Hence

$$\begin{aligned}
\|\Delta_1 \beta^{(t)}\| &= \left\| V^{(t)} \prod_{i=0}^{t-1} \left(I - \Psi W^{(i)} \right) \left(\beta^{(0)} - \hat{\beta}_{2\text{SLS}} \right) \right\| \\
&\leq \|V^{(t)}\| \left\| \prod_{i=0}^{T_0-1} \left(I - \Psi W^{(i)} \right) \right\| \left\| \prod_{i=T_0}^{t-1} \left(I - \Psi W^{(i)} \right) \right\| \|\beta^{(0)} - \hat{\beta}_{2\text{SLS}}\| \\
&< \gamma(\alpha)^t \left\| \prod_{i=0}^{T_0-1} \left(I - \Psi W^{(i)} \right) \right\| (\sqrt{p} + 1) \|\beta^{(0)} - \hat{\beta}_{2\text{SLS}}\| \\
&= \mathcal{O}(\gamma(\alpha)^t).
\end{aligned} \quad (47)$$

Next we consider $\Delta_2 \beta^{(t)}$:

$$\begin{aligned}
\Delta_2 \beta^{(t)} &= \sum_{i=0}^{t-1} \alpha \left[\prod_{j=i+1}^{t-1} \left(I - \alpha \Theta^{(j)\top} Z^\top Z \Theta^{(j)} \right) \right] \Theta^{(i)\top} Z^\top \left(Y - Z \Theta^{(i)} \hat{\beta}_{2\text{SLS}} \right) \\
&= \sum_{i=0}^{t-1} \alpha \left[\prod_{j=i+1}^{t-1} \left(I - \alpha \hat{\Theta}^\top Z^\top Z \hat{\Theta} - \alpha W^{(j)} \right) \right] \left(R^{(i)} + \hat{\Theta} \right)^\top Z^\top \left[Y - Z \left(R^{(i)} + \hat{\Theta} \right) \hat{\beta}_{2\text{SLS}} \right] \\
&= \sum_{i=0}^{t-1} \alpha \left(I - \alpha \hat{\Theta}^\top Z^\top Z \hat{\Theta} \right)^{t-1-i} \prod_{j=i+1}^{t-1} \left[I - \alpha \left(I - \alpha \hat{\Theta}^\top Z^\top Z \hat{\Theta} \right)^{-1} W^{(j)} \right] \\
&\quad \cdot \left(R^{(i)} + \hat{\Theta} \right)^\top Z^\top \left[Y - Z \left(R^{(i)} + \hat{\Theta} \right) \hat{\beta}_{2\text{SLS}} \right] \\
&= \sum_{i=0}^{t-1} \alpha \left[V^{(t-1-i)} \prod_{j=i+1}^{t-1} \left(I - \Psi W^{(j)} \right) \right] \left(R^{(i)} + \hat{\Theta} \right)^\top Z^\top \left[Y - Z \left(R^{(i)} + \hat{\Theta} \right) \hat{\beta}_{2\text{SLS}} \right].
\end{aligned}$$

For convenience, let $\Delta_2 \beta^{(t)} := \Delta_{21} \beta^{(t)} + \Delta_{22} \beta^{(t)}$, where

$$\Delta_{21} \beta^{(t)} := \sum_{i=0}^{t-1} \alpha \left[V^{(t-1-i)} \prod_{j=i+1}^{t-1} \left(I - \Psi W^{(j)} \right) \right] R^{(i)\top} Z^\top \left[Y - Z \left(R^{(i)} + \hat{\Theta} \right) \hat{\beta}_{2\text{SLS}} \right],$$

$$\Delta_{22}\boldsymbol{\beta}^{(t)} := \sum_{i=0}^{t-1} \alpha \left[\mathbf{V}^{(t-1-i)} \prod_{j=i+1}^{t-1} (\mathbf{I} - \boldsymbol{\Psi}\mathbf{W}^{(j)}) \right] \hat{\boldsymbol{\Theta}}^\top \mathbf{Z}^\top \left[\mathbf{Y} - \mathbf{Z} (\mathbf{R}^{(i)} + \hat{\boldsymbol{\Theta}}) \hat{\boldsymbol{\beta}}_{2\text{SLS}} \right].$$

Suppose \tilde{M}_1, \tilde{M}_2 are the upper bounds such that

$$\left\| \mathbf{Z}^\top \left[\mathbf{Y} - \mathbf{Z} (\mathbf{R}^{(i)} + \hat{\boldsymbol{\Theta}}) \hat{\boldsymbol{\beta}}_{2\text{SLS}} \right] \right\| \leq \tilde{M}_1, \quad \forall i = 0, \dots, t-1,$$

$$\left\| \prod_{j=i+1}^{t-1} (\mathbf{I} - \boldsymbol{\Psi}\mathbf{W}^{(j)}) \right\| \leq \tilde{M}_2, \quad \forall i = 0, \dots, t-1.$$

We know such \tilde{M}_1, \tilde{M}_2 exist because of the bounds given by (45) and (46). Let $\tilde{M} = \tilde{M}_1 \tilde{M}_2$. Then

$$\begin{aligned} \left\| \Delta_{21}\boldsymbol{\beta}^{(t)} \right\| &\leq \tilde{M} \left\| \sum_{i=0}^{t-1} \alpha \mathbf{V}^{(t-1-i)} \mathbf{R}^{(i)\top} \right\| \\ &\leq \tilde{M} \alpha \sum_{i=0}^{t-1} \left\| \mathbf{V}^{(t-1-i)} \right\| \left\| \mathbf{R}^{(i)} \right\| \\ &\leq \tilde{M} \alpha \left\| \boldsymbol{\Theta}^{(0)} - \hat{\boldsymbol{\Theta}} \right\| \sum_{i=0}^{t-1} \gamma(\alpha)^{t-1-i} \kappa(\eta)^i \\ &= \tilde{M} \alpha \left\| \boldsymbol{\Theta}^{(0)} - \hat{\boldsymbol{\Theta}} \right\| \sum_{i=0}^{t-1} \gamma(\alpha)^{t-1} \left(\frac{\kappa(\eta)}{\gamma(\alpha)} \right)^i \\ &= \mathcal{O} \left(\frac{\gamma(\alpha)^t - \kappa(\eta)^t}{\gamma(\alpha) - \kappa(\eta)} \right) \\ &\leq \mathcal{O}(\max\{\gamma(\alpha)^t, \kappa(\eta)^t\}), \end{aligned}$$

and similarly,

$$\begin{aligned} \left\| \Delta_{22}\boldsymbol{\beta}^{(t)} \right\| &\leq \tilde{M} \left\| \sum_{i=0}^{t-1} \alpha \mathbf{V}^{(t-1-i)} \hat{\boldsymbol{\Theta}}^\top \right\| \\ &\leq \tilde{M} \alpha \left\| \hat{\boldsymbol{\Theta}} \right\| \sum_{i=0}^{t-1} \left\| \mathbf{V}^{(t-1-i)} \right\| \\ &\leq \tilde{M} \alpha \left\| \hat{\boldsymbol{\Theta}} \right\| \sum_{i=0}^{t-1} \gamma(\alpha)^{t-1-i} \\ &= \mathcal{O}(\gamma(\alpha)^t). \end{aligned}$$

Thus

$$\begin{aligned} \left\| \Delta_2\boldsymbol{\beta}^{(t)} \right\| &= \left\| \Delta_{21}\boldsymbol{\beta}^{(t)} + \Delta_{22}\boldsymbol{\beta}^{(t)} \right\| \\ &\leq \left\| \Delta_{21}\boldsymbol{\beta}^{(t)} \right\| + \left\| \Delta_{22}\boldsymbol{\beta}^{(t)} \right\| \\ &\leq \mathcal{O}(\max\{\gamma(\alpha)^t, \kappa(\eta)^t\}). \end{aligned} \tag{48}$$

Therefore, plugging (47) and (48) into (44), we have:

$$\left\| \boldsymbol{\beta}^{(t)} - \hat{\boldsymbol{\beta}}_{2\text{SLS}} \right\| \leq \mathcal{O}(\max\{\gamma(\alpha)^t, \kappa(\eta)^t\}),$$

Hence completes the proof. \square

B.2 PROOF OF THEOREM 3.2

Proof of Theorem 3.2. For ease of notations, we ignore l in the following proof. Consider the input matrix taking the form:

$$\mathbf{H}^{(0)} = \begin{bmatrix} \mathbf{z}_1 & \cdots & \mathbf{z}_n & \mathbf{z}_{n+1} \\ \mathbf{x}_1 & \cdots & \mathbf{x}_n & \mathbf{x}_{n+1} \\ y_1 & \cdots & y_n & 0 \\ \Theta_{:,1}^{(0)} & \cdots & \Theta_{:,1}^{(0)} & \Theta_{:,1}^{(0)} \\ \vdots & \vdots & \vdots & \vdots \\ \Theta_{:,p}^{(0)} & \cdots & \Theta_{:,p}^{(0)} & \Theta_{:,p}^{(0)} \\ \beta^{(0)} & \cdots & \beta^{(0)} & \beta^{(0)} \\ \hat{\mathbf{x}}_1^{(0)} & \cdots & \hat{\mathbf{x}}_n^{(0)} & \hat{\mathbf{x}}_{n+1}^{(0)} \\ 1 & \cdots & 1 & 1 \\ 1 & \cdots & 1 & 0 \end{bmatrix} \in \mathbb{R}^{D \times (n+1)},$$

i.e., element-wise,

$$\mathbf{h}_i^{(0)} = \left(z_i, \mathbf{x}_i, y_i t_i, \Theta_{:,1}^{(0)}, \dots, \Theta_{:,p}^{(0)}, \beta^{(0)}, \hat{\mathbf{x}}_i^{(0)}, 1, t_i \right)^\top, \quad i = 1, \dots, n+1,$$

where $D = qp + 3p + q + 3$, $t_i := \mathbb{1}\{i \leq n\}$ is the indicator for training sample. We can take any initialization for $\Theta^{(0)}$, $\beta^{(0)}$ and $\hat{\mathbf{x}}^{(0)}$. To avoid abuse of notations, we omit the superscript of those parameters to be updated in the following proof.

Recall the definitions (4) and (5). Our goal is to show that there exists a series of attention parameters $\theta_{\text{ATTN}}^{(1:2)} = \{(\mathbf{Q}_m^{(1:2)}, \mathbf{K}_m^{(1:2)}, \mathbf{V}_m^{(1:2)})\}_{m \in [M]} \subset \mathbb{R}^{D \times D}$ such that $\theta_{\text{ATTN}}^{(1:2)}$ updates Θ, β on the corresponding rows. i.e, if we denote $D_0 := q + p + 1$, the updates on row $D_0 + 1$ to row $D_0 + qp$ correspond to Θ , and the updates on row $D_0 + qp + 1$ to row $D_0 + qp + p$ correspond to β .

1) In the first layer, the transformer updates the current first-stage estimate $\hat{\mathbf{x}}$.

For $m = 2k - 1, k = 1, \dots, p$, define $\mathbf{Q}_m^{(1)}, \mathbf{K}_m^{(1)}, \mathbf{V}_m^{(1)}$ such that:

$$\mathbf{Q}_m^{(1)} \mathbf{h}_i^{(0)} = \begin{bmatrix} z_{i1} \\ \vdots \\ z_{iq} \\ \hat{x}_{ik}^{(0)} \\ \mathbf{0} \end{bmatrix}, \mathbf{K}_m^{(1)} \mathbf{h}_j^{(0)} = \begin{bmatrix} \Theta_{1k}^{(0)} \\ \vdots \\ \Theta_{qk}^{(0)} \\ -1 \\ \mathbf{0} \end{bmatrix}, \mathbf{V}_m^{(1)} \mathbf{h}_j^{(0)} = \mathbf{e}_{D_0 + qp + p + k}. \quad (49)$$

For $m = 2k, k = 1, \dots, p$, define $\mathbf{Q}_m^{(1)}, \mathbf{K}_m^{(1)}, \mathbf{V}_m^{(1)}$ such that:

$$\mathbf{Q}_m^{(1)} \mathbf{h}_i^{(0)} = \begin{bmatrix} -z_{i1} \\ \vdots \\ -z_{iq} \\ \hat{x}_{ik}^{(0)} \\ \mathbf{0} \end{bmatrix}, \mathbf{K}_m^{(1)} \mathbf{h}_j^{(0)} = \begin{bmatrix} \Theta_{1k}^{(0)} \\ \vdots \\ \Theta_{qk}^{(0)} \\ 1 \\ \mathbf{0} \end{bmatrix}, \mathbf{V}_m^{(1)} \mathbf{h}_j^{(0)} = -\mathbf{e}_{D_0 + qp + p + k}, \quad (50)$$

where $\mathbf{e}_j \in \mathbb{R}^D$ is the standard unit vector with only one 1 at the j -th coordinate. Note that the above are just linear transformations on \mathbf{h}_i or \mathbf{h}_j , hence such matrices $\mathbf{Q}_m^{(1)}, \mathbf{K}_m^{(1)}, \mathbf{V}_m^{(1)}$ must exist.

Then we have:

$$\begin{aligned}
\mathbf{h}_i^{(1)} &= \mathbf{h}_i^{(0)} + \sum_{m=1}^{2p} \frac{1}{n+1} \sum_{j=1}^{n+1} \sigma \left(\langle \mathbf{Q}_m^{(1)} \mathbf{h}_i^{(0)}, \mathbf{K}_m^{(1)} \mathbf{h}_j^{(0)} \rangle \right) \cdot \mathbf{V}_m^{(1)} \mathbf{h}_j^{(0)} \\
&= \mathbf{h}_i^{(0)} + \sum_{k=1}^p \frac{1}{n+1} \sum_{j=1}^{n+1} \left[\sigma \left(\sum_{l=1}^q z_{il} \Theta_{lk}^{(0)} - \hat{x}_{ik}^{(0)} \right) - \sigma \left(- \sum_{l=1}^q z_{il} \Theta_{lk}^{(0)} + \hat{x}_{ik}^{(0)} \right) \right] \cdot \mathbf{e}_{D_0+qp+p+k} \\
&= \mathbf{h}_i^{(0)} + \sum_{k=1}^p \left[\sum_{l=1}^q z_{il} \Theta_{lk}^{(0)} - \hat{x}_{ik}^{(0)} \right] \mathbf{e}_{D_0+qp+p+k} \\
&= \mathbf{h}_i^{(0)} + \sum_{k=1}^p \left(\hat{x}_{ik}^{(1)} - \hat{x}_{ik}^{(0)} \right) \mathbf{e}_{D_0+qp+p+k}.
\end{aligned}$$

Thus this layer correctly updates the first-stage prediction values $\hat{x}_1^{(1)}, \dots, \hat{x}_{n+1}^{(1)}$, where $\hat{x}_i^{(1)} := [\mathbf{Z} \Theta^{(0)}]_{i,:} = \sum_{l=1}^q z_{il} \Theta_{l,:}^{(0)}$. We will have:

$$\mathbf{H}^{(1)} = \begin{bmatrix} \mathbf{z}_1 & \cdots & \mathbf{z}_n & \mathbf{z}_{n+1} \\ \mathbf{x}_1 & \cdots & \mathbf{x}_n & \mathbf{x}_{n+1} \\ y_1 & \cdots & y_n & 0 \\ \Theta_{:,1}^{(0)} & \cdots & \Theta_{:,1}^{(0)} & \Theta_{:,1}^{(0)} \\ \vdots & \vdots & \vdots & \vdots \\ \Theta_{:,p}^{(0)} & \cdots & \Theta_{:,p}^{(0)} & \Theta_{:,p}^{(0)} \\ \beta^{(0)} & \cdots & \beta^{(0)} & \beta^{(0)} \\ \hat{\mathbf{x}}_1^{(1)} & \cdots & \hat{\mathbf{x}}_n^{(1)} & \hat{\mathbf{x}}_{n+1}^{(1)} \\ 1 & \cdots & 1 & 1 \\ 1 & \cdots & 1 & 0 \end{bmatrix}.$$

2) In the second layer, the transformer does the gradient updates on the parameters Θ and β .

For $m = 2k - 1, k = 1, \dots, p$, define $\mathbf{Q}_m^{(2)}, \mathbf{K}_m^{(2)}, \mathbf{V}_m^{(2)}$ such that:

$$\mathbf{Q}_m^{(2)} \mathbf{h}_i^{(1)} = \begin{bmatrix} \Theta_{:,k}^{(0)} \\ -1 \\ \vdots \\ \mathbf{0} \end{bmatrix}, \mathbf{K}_m^{(2)} \mathbf{h}_j^{(1)} = \begin{bmatrix} \mathbf{z}_j \\ x_{jk} t_j \\ R(1-t_j) \\ \vdots \\ \mathbf{0} \end{bmatrix}, \mathbf{V}_m^{(2)} \mathbf{h}_j^{(1)} = -(n+1)\eta \sum_{l=1}^q z_{jl} \mathbf{e}_{D_0+(k-1)q+l}. \quad (51)$$

For $m = 2k, k = 1, \dots, p$, define $\mathbf{Q}_m^{(2)}, \mathbf{K}_m^{(2)}, \mathbf{V}_m^{(2)}$ such that:

$$\mathbf{Q}_m^{(2)} \mathbf{h}_i^{(1)} = \begin{bmatrix} -\Theta_{:,k}^{(0)} \\ 1 \\ -1 \\ \vdots \\ \mathbf{0} \end{bmatrix}, \mathbf{K}_m^{(2)} \mathbf{h}_j^{(1)} = \begin{bmatrix} \mathbf{z}_j \\ x_{jk} t_j \\ R(1-t_j) \\ \vdots \\ \mathbf{0} \end{bmatrix}, \mathbf{V}_m^{(2)} \mathbf{h}_j^{(1)} = (n+1)\eta \sum_{l=1}^q z_{jl} \mathbf{e}_{D_0+(k-1)q+l}, \quad (52)$$

where $R = \max_{i=1, \dots, n+1} \{ \|\Theta^{(t)\top} \mathbf{z}_i\| \}$. Then we have:

$$\begin{aligned}
\sigma \left(\langle \mathbf{Q}_{2k-1}^{(2)} \mathbf{h}_i^{(1)}, \mathbf{K}_{2k-1}^{(2)} \mathbf{h}_j^{(1)} \rangle \right) &= \sigma \left(\Theta_{:,k}^{(0)\top} \mathbf{z}_j - x_{jk} t_j - R(1-t_j) \right) \\
&= \sigma \left(\Theta_{:,k}^{(0)\top} \mathbf{z}_j - x_{jk} \right) \mathbb{1}\{t_j = 1\} \\
&= \sigma \left(\Theta_{:,k}^{(0)\top} \mathbf{z}_j - x_{jk} \right) t_j,
\end{aligned}$$

and

$$\begin{aligned}\sigma\left(\langle \mathbf{Q}_{2k}^{(2)} \mathbf{h}_i^{(1)}, \mathbf{K}_{2k}^{(2)} \mathbf{h}_j^{(1)} \rangle\right) &= \sigma\left(-\boldsymbol{\Theta}_{:,k}^{(0)\top} \mathbf{z}_j + x_{jk} t_j - R(1-t_j)\right) \\ &= \sigma\left(-\boldsymbol{\Theta}_{:,k}^{(0)\top} \mathbf{z}_j + x_{jk}\right) \mathbb{1}\{t_j = 1\} \\ &= \sigma\left(-\boldsymbol{\Theta}_{:,k}^{(0)\top} \mathbf{z}_j + x_{jk}\right) t_j.\end{aligned}$$

So that

$$\begin{aligned}& \sum_{m=1}^{2p} \sigma\left(\langle \mathbf{Q}_m^{(2)} \mathbf{h}_i^{(1)}, \mathbf{K}_m^{(2)} \mathbf{h}_j^{(1)} \rangle\right) \mathbf{V}_m^{(2)} \mathbf{h}_j^{(1)} \\ &= -(n+1)t_j \eta \sum_{k=1}^p \left[\sigma\left(\boldsymbol{\Theta}_{:,k}^{(0)\top} \mathbf{z}_j - x_{jk}\right) - \sigma\left(-\boldsymbol{\Theta}_{:,k}^{(0)\top} \mathbf{z}_j + x_{jk}\right) \right] \cdot \sum_{l=1}^q z_{jl} \mathbf{e}_{D_0+(k-1)q+l} \\ &= -(n+1)t_j \eta \sum_{k=1}^p \sum_{l=1}^q z_{jl} \left(\boldsymbol{\Theta}_{:,k}^{(0)\top} \mathbf{z}_j - x_{jk}\right) \mathbf{e}_{D_0+(k-1)q+l}.\end{aligned}$$

Similarly, for $m = 2p+1, 2p+2$, define $\mathbf{Q}_m^{(2)}, \mathbf{K}_m^{(2)}, \mathbf{V}_m^{(2)}$ such that:

$$\mathbf{Q}_{2p+1}^{(2)} \mathbf{h}_i^{(1)} = \begin{bmatrix} \boldsymbol{\beta}^{(0)} \\ -1 \\ -1 \\ \vdots \\ \mathbf{0} \end{bmatrix}, \mathbf{K}_{2p+1}^{(2)} \mathbf{h}_j^{(1)} = \begin{bmatrix} \hat{\mathbf{x}}_j^{(1)} \\ y_j t_j \\ R'(1-t_j) \\ \vdots \\ \mathbf{0} \end{bmatrix}, \mathbf{V}_{2p+1}^{(2)} \mathbf{h}_j^{(1)} = -(n+1)\alpha \sum_{l=1}^p \hat{\mathbf{x}}_{jl}^{(1)} \mathbf{e}_{D_0+qp+l}, \quad (53)$$

$$\mathbf{Q}_{2p+2}^{(2)} \mathbf{h}_i^{(1)} = \begin{bmatrix} -\boldsymbol{\beta}^{(0)} \\ 1 \\ -1 \\ \vdots \\ \mathbf{0} \end{bmatrix}, \mathbf{K}_{2p+2}^{(2)} \mathbf{h}_j^{(1)} = \begin{bmatrix} \hat{\mathbf{x}}_j^{(1)} \\ y_j t_j \\ R'(1-t_j) \\ \vdots \\ \mathbf{0} \end{bmatrix}, \mathbf{V}_{2p+2}^{(2)} \mathbf{h}_j^{(1)} = (n+1)\alpha \sum_{l=1}^p \hat{\mathbf{x}}_{jl}^{(1)} \mathbf{e}_{D_0+qp+l}, \quad (54)$$

where $R' = \max_{i=1, \dots, n+1} \{|\boldsymbol{\beta}^{(t)\top} \mathbf{x}_i|\}$. Then

$$\begin{aligned}\sigma\left(\langle \mathbf{Q}_{2p+1}^{(2)} \mathbf{h}_i^{(1)}, \mathbf{K}_{2p+1}^{(2)} \mathbf{h}_j^{(1)} \rangle\right) &= \sigma\left(\boldsymbol{\beta}^{(0)\top} \hat{\mathbf{x}}_j^{(1)} - y_j t_j - R'(1-t_j)\right) \\ &= \sigma\left(\boldsymbol{\beta}^{(0)\top} \hat{\mathbf{x}}_j^{(1)} - y_j\right) \mathbb{1}\{t_j = 1\} \\ &= \sigma\left(\boldsymbol{\beta}^{(0)\top} \hat{\mathbf{x}}_j^{(1)} - y_j\right) t_j,\end{aligned}$$

and

$$\begin{aligned}\sigma\left(\langle \mathbf{Q}_{2p+2}^{(2)} \mathbf{h}_i^{(1)}, \mathbf{K}_{2p+2}^{(2)} \mathbf{h}_j^{(1)} \rangle\right) &= \sigma\left(-\boldsymbol{\beta}^{(0)\top} \hat{\mathbf{x}}_j^{(1)} + y_j t_j - R'(1-t_j)\right) \\ &= \sigma\left(-\boldsymbol{\beta}^{(0)\top} \hat{\mathbf{x}}_j^{(1)} + y_j\right) \mathbb{1}\{t_j = 1\} \\ &= \sigma\left(-\boldsymbol{\beta}^{(0)\top} \hat{\mathbf{x}}_j^{(1)} + y_j\right) t_j.\end{aligned}$$

So that

$$\begin{aligned}& \sum_{m=2p+1}^{2p+2} \sigma\left(\langle \mathbf{Q}_m^{(2)} \mathbf{h}_i^{(1)}, \mathbf{K}_m^{(2)} \mathbf{h}_j^{(1)} \rangle\right) \mathbf{V}_m^{(2)} \mathbf{h}_j^{(1)} \\ &= -(n+1)t_j \alpha \left[\sigma\left(\boldsymbol{\beta}^{(0)\top} \hat{\mathbf{x}}_j^{(1)} - y_j\right) - \sigma\left(-\boldsymbol{\beta}^{(0)\top} \hat{\mathbf{x}}_j^{(1)} + y_j\right) \right] \cdot \sum_{l=1}^p \hat{\mathbf{x}}_{jl}^{(1)} \mathbf{e}_{D_0+qp+l} \\ &= -(n+1)t_j \alpha \sum_{l=1}^p \hat{\mathbf{x}}_{jl}^{(1)} \left(\boldsymbol{\beta}^{(0)\top} \hat{\mathbf{x}}_j^{(1)} - y_j\right) \mathbf{e}_{D_0+qp+l}.\end{aligned}$$

Thus the final output, for $i = 1, \dots, n + 1$:

$$\begin{aligned}
\mathbf{h}_i^{(2)} &= \mathbf{h}_i^{(1)} + \sum_{m=1}^{2p+2} \frac{1}{n+1} \sum_{j=1}^n \sigma \left(\langle \mathbf{Q}_m^{(2)} \mathbf{h}_i^{(1)}, \mathbf{K}_m^{(2)} \mathbf{h}_j^{(1)} \rangle \right) \mathbf{V}_m^{(2)} \mathbf{h}_j^{(1)} \\
&= \begin{bmatrix} \mathbf{z}_i \\ \mathbf{x}_i \\ y_i t_i \\ \Theta_{:,1}^{(0)} - \eta \sum_{j=1}^n \mathbf{z}_j \left(\mathbf{z}_j^\top \Theta_{:,1}^{(0)} - x_{j1} \right) \\ \vdots \\ \Theta_{:,p}^{(0)} - \eta \sum_{j=1}^n \mathbf{z}_j \left(\mathbf{z}_j^\top \Theta_{:,p}^{(0)} - x_{jp} \right) \\ \beta^{(0)} - \alpha \sum_{j=1}^n \hat{\mathbf{x}}_j^{(1)} \left(\hat{\mathbf{x}}_j^{(1)\top} \beta^{(0)} - y_j \right) \\ \hat{\mathbf{x}}_i^{(1)} \\ 1 \\ t_i \end{bmatrix} \\
&= \begin{bmatrix} \mathbf{z}_i \\ \mathbf{x}_i \\ y_i t_i \\ \Theta_{:,1}^{(0)} - \eta \mathbf{Z}^\top [\mathbf{Z} \Theta^{(0)} - \mathbf{X}]_{:,1} \\ \vdots \\ \Theta_{:,p}^{(0)} - \eta \mathbf{Z}^\top [\mathbf{Z} \Theta^{(0)} - \mathbf{X}]_{:,p} \\ \beta^{(0)} - \alpha \Theta^{(0)\top} \mathbf{Z}^\top (\mathbf{Z} \Theta^{(0)} \beta^{(0)} - \mathbf{y}) \\ \hat{\mathbf{x}}_i^{(1)} \\ 1 \\ t_i \end{bmatrix}.
\end{aligned}$$

This corresponds to a one-step 2SLS GD update of Θ, β , according to (7). Therefore, the final output matrix is:

$$\mathbf{H}^{(2)} = \begin{bmatrix} \mathbf{z}_1 & \cdots & \mathbf{z}_n & \mathbf{z}_{n+1} \\ \mathbf{x}_1 & \cdots & \mathbf{x}_n & \mathbf{x}_{n+1} \\ y_1 & \cdots & y_n & 0 \\ \Theta_{:,1}^{(1)} & \cdots & \Theta_{:,1}^{(1)} & \Theta_{:,1}^{(1)} \\ \vdots & \vdots & \vdots & \vdots \\ \Theta_{:,p}^{(1)} & \cdots & \Theta_{:,p}^{(1)} & \Theta_{:,p}^{(1)} \\ \beta^{(1)} & \cdots & \beta^{(1)} & \beta^{(1)} \\ \hat{\mathbf{x}}_1^{(1)} & \cdots & \hat{\mathbf{x}}_n^{(1)} & \hat{\mathbf{x}}_{n+1}^{(1)} \\ 1 & \cdots & 1 & 1 \\ 1 & \cdots & 1 & 0 \end{bmatrix}.$$

Thus the proof is complete. We further note that in construction steps (49)(50)(51)(52)(53)(54), regardless of the initial values of $\Theta^{(0)}, \beta^{(0)}$, and $\hat{\mathbf{x}}^{(0)}$, the matrices $\mathbf{Q}_m^{(1:2)}, \mathbf{K}_m^{(1:2)}, \mathbf{V}_m^{(1:2)}$ do the same linear transformations on the input vectors. Therefore they are identical across different layers. \square

B.3 PROOF OF THEOREM 3.3

Lemma B.2 (Generalization of pretraining, from Theorem 20 in Bai et al. (2023)). Given optimization problem (12), with probability at least $1 - \zeta$, the solution $\hat{\theta}$ satisfies:

$$L_{\text{ICL}}(\hat{\theta}) \leq \inf_{\theta \in \Theta} L_{\text{ICL}}(\theta) + \mathcal{O} \left(B_y^2 \sqrt{\frac{L^2 (MD^2 + DD') \log(2 + \max\{B_\theta, R, B_y\}) + \log(1/\zeta)}{N}} \right).$$

Proof of Theorem 3.3. We begin by showing the (clipped) 2SLS predictor achieves small excess loss under in-context distribution \mathcal{P} :

$$\begin{aligned}
& \mathbb{E}_{\mathcal{P}} \left[\left(\langle \text{clip}_{B_{\beta}}(\hat{\beta}_{2\text{SLS}}), \mathbf{x}_{n+1} \rangle - y_{n+1} \right)^2 \right] \\
&= \mathbb{E}_{\mathcal{P}} \left[\left(\langle \text{clip}_{B_{\beta}}(\hat{\beta}_{2\text{SLS}}) - \beta, \mathbf{x}_{n+1} \rangle + \langle \beta, \mathbf{x}_{n+1} \rangle - y_{n+1} \right)^2 \right] \\
&= \mathbb{E}_{\mathcal{P}} \left[\langle \text{clip}_{B_{\beta}}(\hat{\beta}_{2\text{SLS}}) - \beta, \mathbf{x}_{n+1} \rangle^2 \right] + 2\mathbb{E}_{\mathcal{P}} \left[\langle \text{clip}_{B_{\beta}}(\hat{\beta}_{2\text{SLS}}) - \beta, \mathbf{x}_{n+1} \rangle (\langle \beta, \mathbf{x}_{n+1} \rangle - y_{n+1}) \right] \\
&\quad + \mathbb{E}_{\mathcal{P}} \left[(\langle \beta, \mathbf{x}_{n+1} \rangle - y_{n+1})^2 \right] \\
&= \underbrace{\mathbb{E}_{\mathcal{P}} \left[\langle \text{clip}_{B_{\beta}}(\hat{\beta}_{2\text{SLS}}) - \beta, \mathbf{x}_{n+1} \rangle^2 \right]}_{\text{Excess Loss}} + \mathbb{E}_{\mathcal{P}} \left[(\langle \beta, \mathbf{x}_{n+1} \rangle - y_{n+1})^2 \right],
\end{aligned}$$

where $\mathbb{E}_{\mathcal{P}}[\langle \text{clip}_{B_{\beta}}(\hat{\beta}_{2\text{SLS}}) - \beta, \mathbf{x}_{n+1} \rangle (\langle \beta, \mathbf{x}_{n+1} \rangle - y_{n+1})] = 0$ follows from the independence between $\langle \text{clip}_{B_{\beta}}(\hat{\beta}_{2\text{SLS}}) - \beta, \mathbf{x}_{n+1} \rangle$ and $(\langle \beta, \mathbf{x}_{n+1} \rangle - y_{n+1})$ with $\mathbb{E}_{\mathcal{P}}[\langle \beta, \mathbf{x}_{n+1} \rangle - y_{n+1}] = \mathbb{E}_{\mathcal{P}}[\epsilon_{n+1}] = 0$.

To bound the excess loss, we have

$$\begin{aligned}
\mathbb{E}_{\mathcal{P}} \left[\langle \text{clip}_{B_{\beta}}(\hat{\beta}_{2\text{SLS}}) - \beta, \mathbf{x}_{n+1} \rangle^2 \right] &= \mathbb{E}_{\mathcal{P}} \left[\left\| \mathbf{x}_{n+1}^{\top} (\text{clip}_{B_{\beta}}(\hat{\beta}_{2\text{SLS}}) - \beta) \right\|^2 \right] \\
&\leq \mathbb{E}_{\mathcal{P}} \left[\|\mathbf{x}_{n+1}\|^2 \left\| \text{clip}_{B_{\beta}}(\hat{\beta}_{2\text{SLS}}) - \beta \right\|^2 \right] \\
&= \mathbb{E}_{\mathcal{P}} \left[\|\mathbf{x}_{n+1}\|^2 \right] \mathbb{E}_{\mathcal{P}} \left[\left\| \text{clip}_{B_{\beta}}(\hat{\beta}_{2\text{SLS}}) - \beta \right\|^2 \right] \\
&\leq \mathcal{O} \left(B_x^2 \left(\frac{q}{n} \left(\frac{B_{\beta}^2}{K} + C^2(n)(\phi^{\top} \Sigma_u \phi + \sigma_{\epsilon}^2) \right) \right) \right),
\end{aligned} \tag{55}$$

where the last inequality follows from (3).

Next, for the ICL loss, we have

$$\begin{aligned}
& L_{\text{ICL}}(\theta) \\
&= \mathbb{E}_{\pi} \mathbb{E}_{\mathcal{P}} \left[\left(\widetilde{\text{TF}}_{\theta}(\mathbf{H}) - y_{n+1} \right)^2 \right] \\
&= \mathbb{E}_{\pi} \mathbb{E}_{\mathcal{P}} \left[\left(\widetilde{\text{TF}}_{\theta}(\mathbf{H}) - \langle \text{clip}_{B_{\beta}}(\hat{\beta}_{2\text{SLS}}), \mathbf{x}_{n+1} \rangle + \langle \text{clip}_{B_{\beta}}(\hat{\beta}_{2\text{SLS}}), \mathbf{x}_{n+1} \rangle - y_{n+1} \right)^2 \right] \\
&= \mathbb{E}_{\pi} \left\{ \mathbb{E}_{\mathcal{P}} \left[\left(\widetilde{\text{TF}}_{\theta}(\mathbf{H}) - \langle \text{clip}_{B_{\beta}}(\hat{\beta}_{2\text{SLS}}), \mathbf{x}_{n+1} \rangle \right)^2 \right] + \mathbb{E}_{\mathcal{P}} \left[\left(\langle \text{clip}_{B_{\beta}}(\hat{\beta}_{2\text{SLS}}), \mathbf{x}_{n+1} \rangle - y_{n+1} \right)^2 \right] \right. \\
&\quad \left. + 2\mathbb{E}_{\mathcal{P}} \left[\left(\widetilde{\text{TF}}_{\theta}(\mathbf{H}) - \langle \text{clip}_{B_{\beta}}(\hat{\beta}_{2\text{SLS}}), \mathbf{x}_{n+1} \rangle \right) \left(\langle \text{clip}_{B_{\beta}}(\hat{\beta}_{2\text{SLS}}), \mathbf{x}_{n+1} \rangle - y_{n+1} \right) \right] \right\} \\
&\leq \mathbb{E}_{\pi} \left\{ \mathbb{E}_{\mathcal{P}} \left[\left(\widetilde{\text{TF}}_{\theta}(\mathbf{H}) - \langle \text{clip}_{B_{\beta}}(\hat{\beta}_{2\text{SLS}}), \mathbf{x}_{n+1} \rangle \right)^2 \right] + \mathbb{E}_{\mathcal{P}} \left[\langle \text{clip}_{B_{\beta}}(\hat{\beta}_{2\text{SLS}}) - \beta, \mathbf{x}_{n+1} \rangle^2 \right] + \mathbb{E}_{\mathcal{P}} \left[(\langle \beta, \mathbf{x}_{n+1} \rangle - y_{n+1})^2 \right] \right. \\
&\quad \left. + 2\mathbb{E}_{\mathcal{P}} \left[\left| \widetilde{\text{TF}}_{\theta}(\mathbf{H}) - \langle \text{clip}_{B_{\beta}}(\hat{\beta}_{2\text{SLS}}), \mathbf{x}_{n+1} \rangle \right| \right] \mathbb{E}_{\mathcal{P}} \left[\left| \langle \text{clip}_{B_{\beta}}(\hat{\beta}_{2\text{SLS}}), \mathbf{x}_{n+1} \rangle - y_{n+1} \right| \right] \right\} \\
&\leq \mathbb{E}_{\pi} \left\{ \mathbb{E}_{\mathcal{P}} \left[\left(\widetilde{\text{TF}}_{\theta}(\mathbf{H}) - \langle \text{clip}_{B_{\beta}}(\hat{\beta}_{2\text{SLS}}), \mathbf{x}_{n+1} \rangle \right)^2 \right] + \mathbb{E}_{\mathcal{P}} \left[\langle \text{clip}_{B_{\beta}}(\hat{\beta}_{2\text{SLS}}) - \beta, \mathbf{x}_{n+1} \rangle^2 \right] + \mathbb{E}_{\mathcal{P}} \left[(\langle \beta, \mathbf{x}_{n+1} \rangle - y_{n+1})^2 \right] \right. \\
&\quad \left. + 2\mathbb{E}_{\mathcal{P}} \left[\left| \widetilde{\text{TF}}_{\theta}(\mathbf{H}) - \langle \text{clip}_{B_{\beta}}(\hat{\beta}_{2\text{SLS}}), \mathbf{x}_{n+1} \rangle \right| \right] \left(\mathbb{E}_{\mathcal{P}} \left[\left| \langle \text{clip}_{B_{\beta}}(\hat{\beta}_{2\text{SLS}}) - \beta, \mathbf{x}_{n+1} \rangle \right| \right] + \mathbb{E}_{\mathcal{P}} \left[\left| \langle \beta, \mathbf{x}_{n+1} \rangle - y_{n+1} \right| \right] \right) \right\}.
\end{aligned} \tag{56}$$

From Corollary 3.1, we know that there exists a $L = 2\bar{L} + 1$ -layer attention-only transformer model θ , with $M = 2(p + 1)$ heads, and embedding dimension $D = qp + 3p + q + 3$, such that for any \mathbf{H} , given any learning rates α, η and Λ as defined in Equation (8), the following holds⁴:

$$\left| \widetilde{\text{TF}}_{\theta}(\mathbf{H}) - \langle \text{clip}_{B_{\beta}}(\hat{\beta}_{2\text{SLS}}), \mathbf{x}_{n+1} \rangle \right| \leq \mathcal{O}\left(B_x \Lambda^{\bar{L}}\right).$$

Denote $\Lambda^* := \min_{\alpha, \eta} \mathbb{E}_{\pi} \mathbb{E}_{\mathcal{P}} [\Lambda | \mathbf{H}, \alpha, \eta]$, then under α^*, η^* , we have:

$$\mathbb{E}_{\pi} \mathbb{E}_{\mathcal{P}} \left[\left| \widetilde{\text{TF}}_{\theta}(\mathbf{H}) - \langle \text{clip}_{B_{\beta}}(\hat{\beta}_{2\text{SLS}}), \mathbf{x}_{n+1} \rangle \right| \right] \leq \mathcal{O}\left(B_x (\Lambda^*)^{\bar{L}}\right), \quad (57)$$

and

$$\mathbb{E}_{\pi} \mathbb{E}_{\mathcal{P}} \left[\left(\widetilde{\text{TF}}_{\theta}(\mathbf{H}) - \langle \text{clip}_{B_{\beta}}(\hat{\beta}_{2\text{SLS}}), \mathbf{x}_{n+1} \rangle \right)^2 \right] \leq \mathcal{O}\left(B_x^2 \mu_{\Lambda, 2}^*\right), \quad (58)$$

where $\mu_{\Lambda, 2}^* := \mathbb{E}_{\pi} \mathbb{E}_{\mathcal{P}} \left[\Lambda^{2\bar{L}} | \mathbf{H}, \alpha^*, \eta^* \right]$ is close to 0.

With condition (13), from Cauchy-Schwarz inequality, we have:

$$\mathbb{E}_{\pi} \mathbb{E}_{\mathcal{P}} \left[\left| \langle \beta, \mathbf{x}_{n+1} \rangle - y_{n+1} \right| \right] \leq \mathbb{E}_{\pi} \left[\sqrt{\mathbb{E}_{\mathcal{P}} (\epsilon_{n+1}^2)} \right] = \mathbb{E}_{\pi} [\sigma_{\epsilon}] \leq \tilde{\sigma}_{\epsilon}. \quad (59)$$

Also, from (55), we have:

$$\mathbb{E}_{\pi} \mathbb{E}_{\mathcal{P}} \left[\langle \text{clip}_{B_{\beta}}(\hat{\beta}_{2\text{SLS}}) - \beta, \mathbf{x}_{n+1} \rangle^2 \right] \leq \mathcal{O}\left(B_x^2 \left(\frac{q}{n} \left(\frac{B_{\beta}^2}{K} + C^2(n) \tilde{\sigma}^2 \right) \right)\right). \quad (60)$$

Further,

$$\begin{aligned} \mathbb{E}_{\pi} \mathbb{E}_{\mathcal{P}} \left[\left| \langle \text{clip}_{B_{\beta}}(\hat{\beta}_{2\text{SLS}}) - \beta, \mathbf{x}_{n+1} \rangle \right| \right] &\leq \sqrt{\mathbb{E}_{\pi} \mathbb{E}_{\mathcal{P}} \left[\langle \text{clip}_{B_{\beta}}(\hat{\beta}_{2\text{SLS}}) - \beta, \mathbf{x}_{n+1} \rangle^2 \right]} \\ &\leq \mathcal{O}\left(B_x \sqrt{\frac{q}{n} \left(\frac{B_{\beta}^2}{K} + C^2(n) \tilde{\sigma}^2 \right)}\right). \end{aligned} \quad (61)$$

Finally, with (57)(58)(59)(60)(61), rearranging the terms in (56), we have:

$$\begin{aligned} &L_{\text{ICL}}(\theta) - \mathbb{E}_{\pi} \mathbb{E}_{\mathcal{P}} \left[(y_{n+1} - \langle \beta, \mathbf{x}_{n+1} \rangle)^2 \right] \\ &\leq \mathcal{O}\left(B_x^2 \left(\mu_{\Lambda, 2}^* + (\Lambda^*)^{\bar{L}} \sqrt{\frac{q}{n} \left(\frac{B_{\beta}^2}{K} + C^2(n) \tilde{\sigma}^2 \right)} + \frac{q}{n} \left(\frac{B_{\beta}^2}{K} + C^2(n) \tilde{\sigma}^2 \right) \right) + B_x (\Lambda^*)^{\bar{L}} \tilde{\sigma}_{\epsilon}\right). \\ &\leq \mathcal{O}\left((\Lambda^*)^{\bar{L}} \left(B_x^2 \sqrt{\frac{q}{n} \left(\frac{B_{\beta}^2}{K} + C^2(n) \tilde{\sigma}^2 \right)} + B_x \tilde{\sigma}_{\epsilon} \right) + B_x^2 \left(\frac{q}{n} \left(\frac{B_{\beta}^2}{K} + C^2(n) \tilde{\sigma}^2 \right) + \mu_{\Lambda, 2}^* \right)\right). \end{aligned} \quad (62)$$

Thus combining Lemma B.2 with (62) completes the proof. \square

C ADDITIONAL EXPERIMENTS AND DISCUSSIONS

For all experiments in this section, to be consistent with our main experiment in Section 4, we generate $n = 50$ training samples with $p = 5, q = 10$, following the data generating process described in Algorithm 1. The task parameters $\Theta, \beta, \Phi, \phi$ are sampled from standard Gaussian distribution, the covariance matrices $\Sigma_z, \Sigma_u, \Sigma_{\omega}$ are set to be identity matrices, and the noise level σ_{ϵ} is set to 1.

⁴The clipping bound on $\hat{\beta}_{2\text{SLS}}$ can be matched by adjusting the clipping threshold on the last layer of $\widetilde{\text{TF}}_{\theta}$.

C.1 SIMULATIONS VERIFYING THEOREM 3.1

We use the GD-based 2SLS method (7) to estimate the causal effect β . For the simulated data, we calculate the following metrics:

$$\frac{2}{\sigma_{\max}^2(\mathbf{Z}\hat{\Theta})} = 0.0016, \quad \frac{2}{\sigma_{\max}^2(\mathbf{Z})} = 0.0212.$$

By Theorem 3.1, the gradient descent converges when $\alpha \in (0, 0.0016)$ and $\eta \in (0, 0.0212)$. The overall convergence rate is determined by $\Lambda := \max\{\gamma(\alpha), \kappa(\eta)\}$, where

$$\begin{aligned} \gamma(\alpha) &:= \rho(\mathbf{I} - \alpha\hat{\Theta}^\top \mathbf{Z}^\top \mathbf{Z}\hat{\Theta}), \\ \kappa(\eta) &:= \rho(\mathbf{I} - \eta\mathbf{Z}^\top \mathbf{Z}). \end{aligned}$$

We first set $\alpha = 0.0012$ and vary η . The corresponding convergence rates are determined by $\Lambda = \max(0.87, \kappa(\eta))$. Next, we set $\eta = 0.01$ and vary α . The corresponding convergence curves are determined by $\Lambda = \max(\gamma(\alpha), 0.82)$. We compare the estimates $\hat{\beta}^{(t)}$ with the 2SLS estimate $\hat{\beta}_{2SLS}$ as the iteration proceeds. The convergence results are shown in Figure 3.

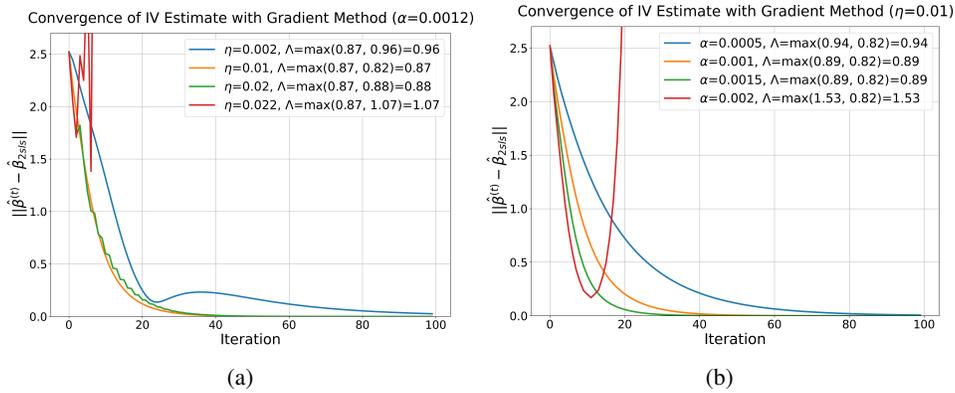


Figure 3: The convergence of the GD-based 2SLS method with (a) fixed $\alpha = 0.0012$ and varying η and (b) fixed $\eta = 0.01$ and varying α .

The results in Figure 3 are consistent with our theoretical analysis in Theorem 3.1. It is worth noting that in Figure 3a, when η is relatively large (or small), the convergence curves exhibit some suiggly patterns. This is due to the innerloop updates (7a) are converging faster (or slower) than the outer loop updates (7b). However, the overall convergence rate is still determined by Λ . This pattern doesn't appear in Figure 3b as we set η to be a moderate value, which ensures that the inner loop and outer loop converge synchronously.

Next, we show the bias of the GD estimator. For better convergence, we set $\alpha^* = \frac{1}{\sigma_{\max}^2(\mathbf{Z}\hat{\Theta})}$ and $\eta^* = \frac{1}{\sigma_{\max}^2(\mathbf{Z})}$. We compare the biases of the GD estimator with $n = 50, 100, 150$ in-context samples. The results are shown in Figure 4.

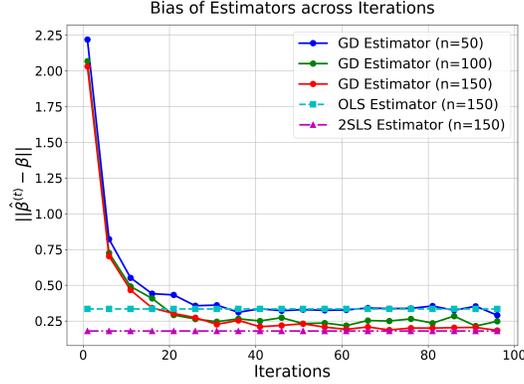


Figure 4: The convergence of the GD-based 2SLS method with $\alpha^* = \frac{1}{\sigma_{\max}^2(\mathbf{Z}\hat{\Theta})}$ and $\eta^* = \frac{1}{\sigma_{\max}^2(\mathbf{Z})}$. The biases of 2SLS estimator and OLS estimator at $n = 150$ are plotted for comparison.

C.2 DISCUSSIONS ON 2SLS WITH ℓ_2 -REGULARIZATION

In this section, we briefly discuss a generalization of our analysis to the case where the 2SLS estimator is regularized by the ℓ_2 penalty (Ridge 2SLS). For this case, the bi-level optimization problem Equation (6) is modified as follows:

$$\min_{\beta} \mathcal{L}(\beta) = \frac{1}{n} \sum_{i=1}^n (y_i - \mathbf{z}_i^\top \hat{\Theta} \beta)^2 + \frac{1}{2} \lambda \|\beta\|^2,$$

$$\text{where } \hat{\Theta} := \arg \min_{\Theta} \frac{1}{n} \sum_{j=1}^n (\mathbf{x}_j - \mathbf{z}_j^\top \Theta)^2 + \frac{1}{2} \tau \|\Theta\|_F^2,$$

where $\lambda, \tau \geq 0$ are regularization parameters. To solve this problem, the GD updates in Equation (7) is modified as follows:

$$\Theta^{(t+1)} = \Theta^{(t)} - \eta \left[\mathbf{Z}^\top (\mathbf{Z} \Theta^{(t)} - \mathbf{X}) + \tau \Theta^{(t)} \right], \quad (63a)$$

$$\beta^{(t+1)} = \beta^{(t)} - \alpha \left[\Theta^{(t)\top} \mathbf{Z}^\top (\mathbf{Z} \Theta^{(t)} \beta^{(t)} - \mathbf{Y}) + \lambda \beta^{(t)} \right]. \quad (63b)$$

The only difference between Equation (7) and Equation (63) is the additional terms $\tau \Theta^{(t)}$ in Equation (63a), and $\lambda \beta^{(t)}$ in Equation (63b). The convergence analysis of the ℓ_2 -regularized GD updates in Equation (63) can be conducted in a similar manner as in Theorem 3.1. The only difference is that the convergence rate Λ is now determined by the spectral radiuses of $\mathbf{I} - \eta(\mathbf{Z}^\top \mathbf{Z} + \tau \mathbf{I})$ and $\mathbf{I} - \alpha(\hat{\Theta}^\top \mathbf{Z}^\top \mathbf{Z} \hat{\Theta} + \lambda \mathbf{I})$, respectively.

With the same configuration as in Theorem 3.2 but adding $p + 1$ attention heads in the second layer (i.e. the second layer needs $3p + 3$ heads), we can show that transformers are able to implement the ℓ_2 -regularized GD updates in Equation (63). The proof follows directly from Appendix B.2, with the construction of the new attention heads in the second layer as follows.

For $m = 2p + 2 + k, k = 1, \dots, p$, define $\mathbf{Q}_m^{(2)}, \mathbf{K}_m^{(2)}, \mathbf{V}_m^{(2)}$ such that:

$$\mathbf{Q}_m^{(2)} \mathbf{h}_i^{(1)} = \begin{bmatrix} 1 \\ \mathbf{0} \end{bmatrix}, \mathbf{K}_m^{(2)} \mathbf{h}_i^{(1)} = \begin{bmatrix} 1 \\ \mathbf{0} \end{bmatrix}, \mathbf{V}_m^{(2)} \mathbf{h}_i^{(1)} = -\eta \tau \sum_{l=1}^q \Theta_{lk}^{(0)} \mathbf{e}_{D_0 + (k-1)q + l}.$$

For $m = 3p + 3$, define $\mathbf{Q}_{3p+3}^{(2)}, \mathbf{K}_{3p+3}^{(2)}, \mathbf{V}_{3p+3}^{(2)}$ such that:

$$\mathbf{Q}_{3p+3}^{(2)} \mathbf{h}_i^{(1)} = \begin{bmatrix} 1 \\ \mathbf{0} \end{bmatrix}, \mathbf{K}_{3p+3}^{(2)} \mathbf{h}_i^{(1)} = \begin{bmatrix} 1 \\ \mathbf{0} \end{bmatrix}, \mathbf{V}_{3p+3}^{(2)} \mathbf{h}_i^{(1)} = -\alpha \lambda \sum_{l=1}^q \beta_l^{(0)} \mathbf{e}_{D_0 + qp + l}.$$

Then the remaining proof follows exactly the same as Appendix B.2. This result indicates that transformers are potentially capable of handling multicollinearity in IV regression problem. We conduct experiments to validate this and the results are shown in Appendix C.3.

C.3 EXPERIMENTS ON MULTICOLLINEARITY IV PROBLEM

As a supplement to Section 4.2, we examine the case where multicollinearity occurs in the IV regression problem. We generate the test prompts in the same way using Algorithm 1, but introduce multicollinearity in the endogenous variable \mathbf{x} and instrument \mathbf{z} .

Specifically, we first generate 4 columns of \mathbf{X} , and 9 columns of \mathbf{Z} , and then set $\mathbf{X}_{:,5} \sim \mathcal{N}(2\mathbf{X}_{:,4}, 10^{-6}\mathbf{I}_n)$, and $\mathbf{Z}_{:,10} \sim \mathcal{N}(2\mathbf{Z}_{:,9}, 10^{-6}\mathbf{I}_n)$. The results are shown in Figure 5a. As shown in Figure 5a, both ordinary OLS and 2SLS estimators fail to estimate the coefficients, while the trained transformer model is still able to provide consistent predictions and coefficient estimates.

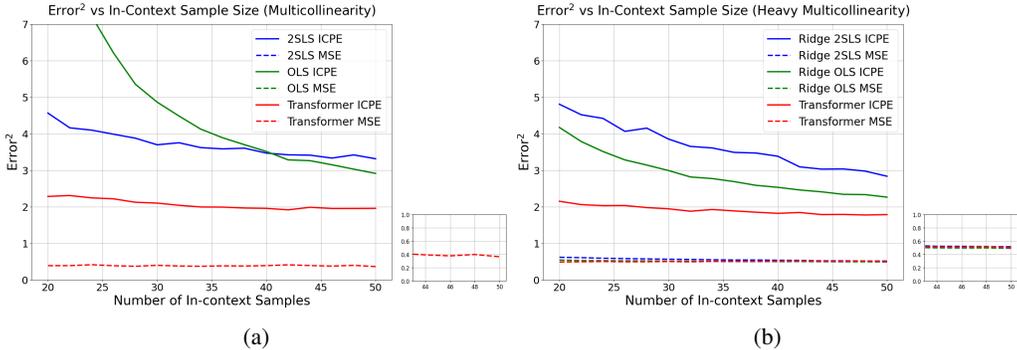


Figure 5: The ICL performance of the trained transformer model in endogeneity tasks with multicollinearity. (a) 1 collinear column in \mathbf{X} , and 1 collinear column in \mathbf{Z} . Note that the coefficient MSEs for 2SLS and OLS are both out of range. (b) 2 collinear columns in \mathbf{X} , and 5 collinear columns in \mathbf{Z} . We compare the performance to the ℓ_2 -regularized 2SLS and OLS estimators. The curves are averaged over 500 simulations.

We further examine a more difficult case where heavy multicollinearity occurs. We first generate 3 column of \mathbf{X} , and 5 column of \mathbf{Z} , and then set $\mathbf{X}_{:,j} \sim \mathcal{N}(2\mathbf{X}_{:,j-2}, 10^{-6}\mathbf{I}_n)$ for $j = 4, 5$, and $\mathbf{Z}_{:,j} \sim \mathcal{N}(2\mathbf{Z}_{:,j-5}, 10^{-6}\mathbf{I}_n)$ for $j = 6, 7, 8, 9, 10$. For better comparisons, we now compare the performance of the trained transformer model to the ℓ_2 -regularized 2SLS and OLS estimators (with all regularization parameters set to 1). The results are shown in Figure 5b.

These results suggest that the trained transformer model is capable to handle multicollinearity in IV regression problems, even though it has not been specifically trained with multicollinearity tasks.

C.4 EXPERIMENTS ON COMPLEX NON-LINEAR IV PROBLEM

As a supplement to Section 4.2, we examine a more complex scenario where the instrument \mathbf{z} has non-linear effect on the endogenous regressor \mathbf{x} . We consider the following data generating process:

$$y = \langle \boldsymbol{\beta}, \mathbf{x} \rangle + \epsilon_1, \quad \text{and} \quad \mathbf{x} = g(\mathbf{z}) + \epsilon_2,$$

where $g : \mathbb{R}^q \rightarrow \mathbb{R}^p$ is a two-layer fully connected neural network with ReLU activation function. Similar to Section 4, the test prompts are generated using Algorithm 1, with all task parameters and weights of neural network sampled from standard Gaussian distribution. The results are shown in Figure 6. From this figure, we can see that the trained transformer model still achieves optimal performance in this complex non-linear setting.

C.5 EXPERIMENTS ON VARYING ENDOGENEITY STRENGTH

As a supplement to Section 4.2, we examine the performance of the trained transformer model in standard IV tasks with varying endogeneity strengths. The strength of endogeneity is determined

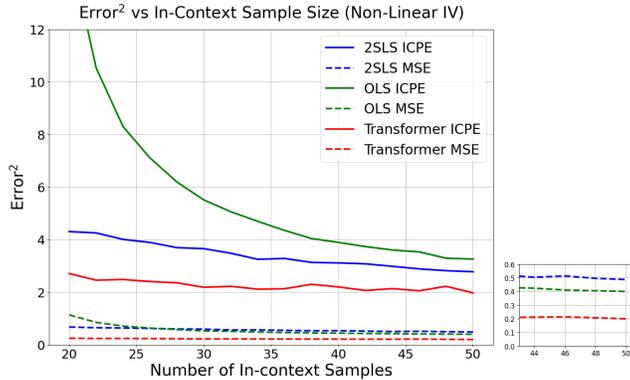


Figure 6: The ICL performance of the trained transformer model in complex non-linear endogeneity tasks where the IV has non-linear effect on the endogenous variable. The curves are averaged over 500 simulations.

by the correlation between x and the endogenous error ϵ_1 . To vary the endogeneity strength, in Algorithm 1, we multiply u by a factor $r \in (0, 2)$ when generating test prompts. The results are shown in Figure 7, which illustrates that the trained transformer model is comparable with the optimal 2SLS estimator in these standard IV tasks, regardless of the endogeneity level.

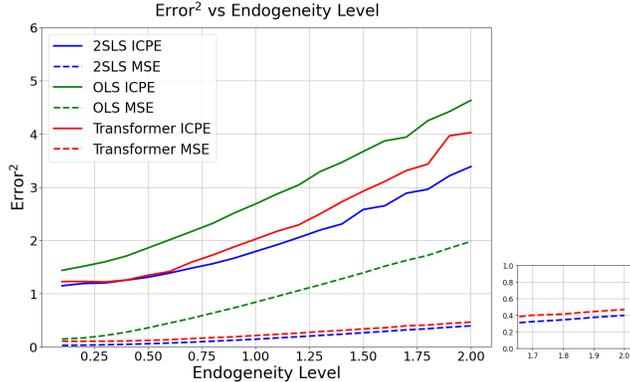


Figure 7: The ICL performance of the trained transformer model in tasks with varying endogeneity strengths. The curves are averaged over 500 simulations.

C.6 EXPERIMENTS ON REAL-WORLD DATASET

In this section we provide an example to illustrate how the pretrained transformer model can be applied to a real-world dataset. We use the dataset from the study of Angrist & Evans (1998). This study investigates the effect of childbearing on labor supply. For demonstration purpose, we consider a simplified setup. We focus on a subset of the dataset that contains 6421 samples from Alabama. The outcome variable y is mother’s labor supply (number of working weeks in a year divided by 52), the endogenous variable x is the number of children (≥ 2), and the instrument z is an indicator variable of whether the first and second children are of the same sex⁵.

⁵Research found that parents of same-sex siblings are significantly more likely to go on to have an additional child (Westoff & Parke, 1972), while it is not directly correlated with mother’s labor supply as mixture of sex of the first two children can be considered as randomly assigned.

For each run we randomly select 50 samples from the dataset, and make the boxplot of the estimated β over 500 runs⁶. As the ground truth effect β is unknown, we compare them to the OLS and 2SLS estimates over all samples. The results are shown in Figure 8.

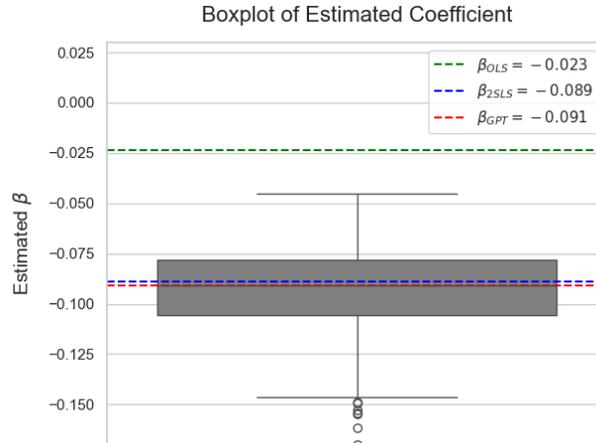


Figure 8: The boxplot of transformer’s estimates over 500 runs on the labor supply dataset, comparing to the OLS and 2SLS estimates. β_{GPT} is taken to be the median of all runs. The gray box represents the interquartile range, where the middle 50% of the estimated values fall. The whiskers of the box indicate the spread of the estimates. Any points falling outside of the whisker can be considered as outliers.

The final estimate $\beta_{GPT} = -0.091$, which suggests that with each increase in the number of children, the mother’s labor supply is expected to drop 9.1% (approximately 4.73 weeks per year). This result is closer to the 2SLS estimate $\beta_{2SLS} = -0.089$ than the OLS estimate $\beta_{OLS} = -0.023$. This example demonstrates the potential of the pretrained transformer model in handling real-world IV problems.

C.7 EXPERIMENTAL DETAIL

The training of the transformer in our experiment was conducted on a Windows 11 machine with the following specifications:

- GPU: NVIDIA GeForce RTX 4090
- CPU: Intel Core i9-14900KF
- Memory: 32 GB DDR5, 5600MHz

The training process took around 10 hours.

⁶For large enough model that can fit in the entire dataset, this step can be ignored. As shown in the simulation study in Section 4, a single estimate is expected to perform at least as good as 2SLS estimator, given the same number of samples.