# Judicial Sentencing Prediction Based on Hybrid Models and Two-Stage Learning Algorithms

Ruifen Dai, Xin Zheng, Fang Wang, and Lei Guo, *Fellow, IEEE*

*Abstract*—The investigation of legal judgment prediction (LJP), such as sentencing prediction, has attracted broad attention for its potential to promote judicial fairness, making the accuracy and reliability of its computation result an increasingly critical concern. In view of this, we present a new sentencing model that shares both legal logic interpretability and strong prediction capability by introducing a two-stage learning algorithm. Specifically, we first construct a hybrid model that synthesizes a mechanism model based on the main factors for sentencing with a neural network modeling possible uncertain features. We then propose a two-stage learning algorithm: First, an adaptive stochastic gradient (ASG) algorithm is used to get good estimates for the unknown parameters in the mechanistic component of the hybrid model. Then, the Adam optimizer tunes all parameters to enhance the predictive performance of the entire hybrid model. The asymptotic convergence of the ASG-based adaptive predictor is established without requiring any excitation data conditions, thereby providing a good initial parameter estimate for prediction. Based on this, the fast-converging Adam optimizer further refines the parameters to enhance overall prediction accuracy. Experiments on a real-world dataset of intentional injury cases in China show that our new hybrid model combined with our two-stage ASG-Adam algorithm, outperforms the existing related methods in sentencing prediction performance, including those based on neural networks and saturated mechanism models.

*Index Terms*—Sentencing Prediction, Saturated Mechanism Model, Neural Network, Two-Stage ASG-Adam Algorithm, Prediction Error Minimization.

## I. Introduction

IN recent years, the investigation of various problems in legal judgment prediction (LJP), such as sentencing prediction, has attracted broad research attention due to its vital role in enhancing fairness in the judicial system ( [2]- [10]). Among these problems, ensuring the high reliability and accuracy of sentencing prediction results is an important requirement in judicial practice. This requirement calls for sentencing models that integrate the sentencing logic for interpretability and capture case-specific uncertainties not specified explicitly in law for prediction capability. However, such models have been rarely developed in the literature. Moreover, most existing theories on the current mainstream gradient-based learning algorithms are not directly applicable to sentencing prediction, since they rely on common yet stringent statistical assumptions on the data, such as the independent and identically distributed (i.i.d.) data assumption, which are not the case for real-world judicial datasets. Therefore, developing sentencing prediction models based on judicial logic and establishing the performance guarantees of the associated learning algorithms, hold significant value in both theory and practice.

Although most existing LJP studies have not adequately addressed the aforementioned challenges, their efforts toward achieving high-precision sentencing prediction offer valuable insights (see, e.g., [1]– [7], [11]– [12], [17], [18]), which can be broadly summarized in terms of model construction and algorithm design in the following:

First, extensive research has been conducted on using machine learning models for sentencing prediction. For example, [1] introduced a convolutional neural network (CNN)-based model for sentencing classification prediction. To address the need to model subtask dependencies, [2] employed a multi-task learning framework TopJudge that represents the dependencies among LJP subtasks using a directed acyclic graph. To further capture the relationships between subtasks, [3] proposed a multi-perspective bi-feedback network model enhanced by a word collocation attention mechanism, which reflects subtask dependencies through a topology-aware design. Subsequently, [4] modeled LJP as a node classification problem over a global consistency graph. [5] presented the LADAN model, which combines a graph neural network (GNN) and an attention mechanism to distinguish confusing law articles for more accurate LJP. In addition, [6] introduced a NeurJudge model that sets BERT (see [13]) as a judgment encoder and leverages intermediate

Ruifen Dai and Fang Wang are with the Data Science Institute, Shandong University, Jinan 250100, China. (e-mails: dairuifen@mail.sdu.edu.cn, wangfang226@sdu.edu.cn).

Xin Zheng and Lei Guo are with the State Key Laboratory of Mathematical Sciences, Academy of Mathematics and Systems Science, Chinese Academy of Sciences, Beijing 100190, China, and also with the School of Mathematical Sciences, University of Chinese Academy of Sciences, Beijing 100049, China. (e-mails: zhengxin2021@amss.ac.cn, lguo@amss.ac.cn).

subtask results to partition fact descriptions into circumstances for guiding subsequent subtask predictions. [7] presented a QAjudge model composed of a question net and an answer net to visualize the prediction processes.

Despite the widespread application of the machine learning sentencing models, they fail to fully incorporate the sentencing logic in the sentencing guidelines. In view of this, [11] proposed an interpretable saturated mechanism model based on the sentencing logic of "starting point—benchmark sentence—announced sentence", and applied it to sentencing prediction for the crime of intentional injury. To improve sentencing prediction accuracy, [12] incorporated saturated boundaries to ensure that neural network–predicted sentences fall within the statutory sentencing range. [31] and [32] refined the sentencing mechanism model proposed in [11] by incorporating the temporal logic underlying sentencing adjustments for benchmark sentences. The impact that not explicitly encoded in sentencing guidelines, although potentially varying across different cases, is coarsely represented by a fixed constant in [31] and [32].

Second, many efforts have been devoted to designing algorithms for sentencing prediction. For instance, [17] proposed a two-step Newton-type adaptive algorithm and analyzed the related performances for both parameter estimation and sentencing prediction, without resorting to the traditional persistence of excitation (PE) on the system data. Subsequently, [18] designed a more robust two-step weighted $l_1$-based Newton-type algorithm to improve prediction performance and established its global convergence. Both algorithms in [17]- [18] were applied to sentencing prediction tasks based on the saturated mechanism model proposed in [11]. However, all the above algorithms incur high computational costs when processing large-scale and high-dimensional judicial data due to their second-order nature. Different from the Newton-type algorithms, the gradient descent algorithms and their variants, such as Adam [15], Adadelta [16], and AdamW [14], are also used for sentencing prediction based on neural network models (e.g., [1]- [7], [12]). However, neural networks are just mathematical approximation and can hardly provide judicial interpretability for the sentencing results.

To overcome the aforementioned shortcomings, we establish a new sentencing model that possesses both high legal interpretability and strong predictive capability by introducing a two-stage gradient algorithm. The main contributions can be summarized as follows:

- We will introduce a hybrid model that integrates the sentencing-logic-based mechanism model with a neural network, which can not only reflect the main factors for sentencing, but also can reflect the influence of possible uncertain factors not explicitly encoded in the sentencing guidelines.

- We will propose a two-stage learning (TSL) algorithm for the prediction of the hybrid sentencing model. To be specific, an adaptive stochastic gradient (ASG) algorithm will provide a good initial value of the mechanistic component parameters for the Adam optimizer in the next stage to fine-tune all parameters to improve the sentencing prediction performance.

- We will rigorously establish the global asymptotic convergence of the ASG-based adaptive predictor without requiring any excitation data conditions. This theoretical guarantee ensures that the ASG-based initialization provides a reliable approximation to globally optimal prediction, thereby offering a solid basis for the refined prediction by the Adam optimizer.

- Empirical experiments using a real-world judgment dataset for the crime of intentional injury (CII) will show that our proposed hybrid model and two-stage algorithm can achieve high accuracy in sentencing prediction, outperforming other known related methods.

The structure of this paper is as follows: Section II will introduce the new sentencing model. Section III will present the two-stage ASG-Adam algorithm. Section IV will provide the global asymptotic convergence theory of the ASG algorithm whose proof will be shown in the appendix. In Section V, we will demonstrate the advantages of the proposed model and algorithm using a real-world sentencing dataset. The concluding remarks will be presented in the final section.

## II. A Hybrid SMNN Model

In this section, we propose a hybrid model that integrates the sentencing mechanism (SM) in Chinese Criminal Law with a neural network (NN) accounting for possible uncertainties in sentencing, abbreviated as SMNN model, as follows:

$$z_{k+1} = S_k \Bigg( \underbrace{\left[ a_k + b x_k^{(1)} + c x_k^{(2)} \right]}_{\text{benchmark sentence}}$$
$$\times \underbrace{\prod_{i=1}^{m_1} (1 + p_i v_k^{(i)})}_{\text{influence of primary factors}} \quad (1)$$
$$\times \underbrace{\left[ 1 + \sum_{j=1}^{m_2} q_j u_k^{(j)} + e_k \right]}_{\text{influence of other factors}} + w_{k+1} \Bigg),$$

where the time-varying bias term $e_k$ and the saturated function $S_k(\cdot)$ are defined respectively as follows:

$$e_k = \Gamma \sigma \left( B \sigma \left( A \eta_k + b^{(1)} \right) + b^{(2)} \right) + b^{(3)}, \quad (2)$$

$$S_k(x) = \begin{cases} L_k, & x < L_k; \\ x, & L_k \le x \le N_k; \\ N_k, & x > N_k, \end{cases} \quad (3)$$

and where $z_k \in \mathbb{R}$ is the pronounced sentence, $a_k \in \mathbb{R}$ is the sentencing starting point, and $x_k^{(i)} \in \mathbb{R}, i = 1, 2$, represent the factors determining penalty amounts, e.g., in the CII [1], they represent the number of seriously and minorly injured victims, respectively. $b, c \in \mathbb{R}$ quantify the additional sentence for the offender corresponding to a one-unit increase in $x_k^{(i)}$. $v_k^{(i)}, i = 1, \cdots, m_1$ denote primary sentencing factors with application priority, and $u_k^{(j)}, j = 1, \cdots, m_2$ denote other sentencing factors. $p_i, i = 1, \cdots, m_1$ and $q_j, j = 1, \cdots, m_2$ are unknown weighting parameters for $v_k^{(i)}$ and $u_k^{(j)}$, respectively. $m_1$ and $m_2$ are the number of primary sentencing factors and other sentencing factors, respectively. $w_k \in \mathbb{R}$ denotes potential random noise effects. $e_k \in \mathbb{R}$ in (2) is the bias term reflecting the comprehensive influence of possible factors not specified explicitly in the law, and $A \in \mathbb{R}^{m \times m_3}, B \in \mathbb{R}^{m \times m}, \Gamma \in \mathbb{R}^{1 \times m}$ are unknown weight parameter matrices or vectors, and $b^{(1)} \in \mathbb{R}^m, b^{(2)} \in \mathbb{R}^m, b^{(3)} \in \mathbb{R}$ are unknown bias parameters in the neural network. Here, $m$ represents the number of neurons in each of the two hidden layers. $m_3$ denotes the number of possible factors that not specified in the law. $\eta_k \in \mathbb{R}^{m_3}$ denotes the vector of factors that not explicitly encoded in the law of the $k$-th case's sentencing. $\sigma(X)$ is an activation function, such as the Rectified Linear Unit (ReLU) function, defined as $(\max(x_1, 0), \max(x_2, 0), \ldots, \max(x_m, 0))^\tau$ with $x_i$ being the element of $X$ in the $i$-th row, $i = 1, \cdots, m$. $L_k \in \mathbb{R}$ and $N_k \in \mathbb{R}$ in (3) are the lower and upper bounds for the announced sentence, respectively, which are prescribed in Chinese Criminal Law and may vary for different crimes and criminal cases.

**Remark 1** *The SMNN model (1) is mainly based on Chinese Criminal Law and sentencing guidelines[2], it specifically follows the sentencing logic of "starting point—benchmark sentence—announced sentence" and incorporates temporal adjustments for benchmark sentences (see [31], [32] for more details). Moreover, the model employs a neural network to approximate the comprehensive influences of factors that are not explicitly encoded in the law. Unlike treating the bias term as a fixed constant in previous studies ( [31] and [32]), the neural network can approximate the case-specific influences, enhancing the model's flexibility and predictive performance.*

*The SMNN model (1) will degenerate to the mechanism model in [11], if we take $v_k^{(i)} = 0, i = 1, \cdots, m_1$ and replace $e_k$ with a constant $e$.*

**Remark 2** *The saturated function (3) ensures that the sentences are restricted to the statutory penalty range, and when a sentence exceeds the upper limit $N_k$ or falls below the lower limit $L_k$, the final judgment is capped at $N_k$ or $L_k$, respectively.*

*Besides, the saturated phenomena described by (3) are also commonly found in various application fields, including engineering systems ( [26]- [27]), economic behavior analysis ( [28]- [29]), biomedical systems ( [30]), and others. These applications collectively emphasize the importance of saturated scenarios in complex nonlinear modeling and analysis.*

**Remark 3** *The SMNN model mainly comprises two parts: (i) a mechanistic component grounded in the sentencing logic, and (ii) a neural network that captures the influence of possible residual factors (RF) not explicitly encoded in sentencing law. Since the mechanistic component already encodes the factors and rules explicitly stated in the sentencing law, which plays a dominating role in sentencing, the influence of the RF may not be significant. Therefore, we first assume the bias term $e_k = e$ to be a constant reflecting the averaged effect of the RF, in order to facilitate a more straightforward estimation of the mechanistic component. After this stage, the highly nonlinear neural network will be introduced to enhance the prediction accuracy for each case.*

To avoid non-convex optimization problems to get better estimators of the mechanistic component parameters, we rewrite the internal structure of (1) as a linearly parameterized regression by increasing the dimension of both regression vectors and parameter vectors. Here the bias term $e_k = e$. To be specific, set the expanded regressor and parameter vector as follows:

$$\phi_k = \left[a_k \phi_{1k}^T, a_k(\phi_{2k} \otimes \phi_{1k})^T, x_k^{(1)} \phi_{1k}^T, \right.$$
$$\left. x_k^{(1)}(\phi_{2k} \otimes \phi_{1k})^T, x_k^{(2)} \phi_{1k}^T, x_k^{(2)}(\phi_{2k} \otimes \phi_{1k})^T\right]^T, \quad (4)$$

$$\theta = \left[(1+e)\vartheta_1^T, (\vartheta_2 \otimes \vartheta_1)^T, b(1+e)\vartheta_1^T, \right.$$
$$\left. b(\vartheta_2 \otimes \vartheta_1)^T, c(1+e)\vartheta_1^T, c(\vartheta_2 \otimes \vartheta_1)^T\right]^T, \quad (5)$$

where

$$\phi_{1k} = [1, z_k^{(1)}, \cdots, z_k^{(m_1)}, z_k^{(1)} z_k^{(2)}, \cdots,$$
$$z_k^{(m_1-1)} z_k^{(m_1)}, \cdots, z_k^{(1)} \cdots z_k^{(m_1)}]^T, \quad (6)$$

$$\phi_{2k} = [u_k^{(1)}, u_k^{(2)}, \cdots, u_k^{(m_2)}]^T, \quad (7)$$

$$\vartheta_1 = [1, p_1, \cdots, p_{m_1}, p_1 p_2, \cdots, p_{m_1-1} p_{m_1},$$
$$\cdots, p_1 \cdots p_{m_1}]^T, \quad (8)$$

$$\vartheta_2 = [q_1, \cdots, q_{m_2}]^T, \quad (9)$$

and the Kronecker product of two vectors is defined as $\mathbf{a} \otimes \mathbf{b} = [a_1\mathbf{b}, \ a_2\mathbf{b}, \ \cdots, \ a_m\mathbf{b}]$. Then (1) can be rewritten as follows:

$$z_{k+1} = S_k(\phi_k^T \theta + w_{k+1}). \tag{10}$$

For the above new parameterized model (10), the dimension of the unknown parameter $\theta$ is $2^{m_1}(3+3m_2)$, which grows exponentially with $m_1$, resulting in a huge demand for computational resources to handle such high-dimensional inputs. Taking the CII example as detailed below, the dimension of unknown parameter vector is 565,248, rendering many Newton-type methods (e.g., [17], [18]) computationally infeasible due to excessive memory and time complexity.

## III. A GRADIENT-BASED TWO-STAGE ALGORITHM

In this section, we propose the two-stage gradient learning algorithm for the SMNN model (1).

To better introduce the algorithm, we first introduce the following notations and assumptions.

### A. Notations and Assumptions

**Notations.** Throughout the sequel, $\|\cdot\|_1$ and $\|\cdot\|$ denote the 1-norm and Euclidean norm of vectors or matrices, respectively. The maximum and minimum eigenvalues of a square matrix $X$ are denoted by $\lambda_{\max}\{X\}$ and $\lambda_{\min}\{X\}$, respectively. Let $\{\mathcal{F}_k\}$ be a non-decreasing sequence of $\sigma$-algebras, along with the associated conditional expectation operator $\mathbb{E}[\cdot \mid \mathcal{F}_k]$. Moreover, for any two sequences $\{a_n\}, \{b_n\}$ with $b_n > 0$, $a_n = O(b_n)$ means that there exists a constant $N > 0$ such that $|a_n|/b_n \leq N$ for all $n > 0$, and $a_n = o(b_n)$ means that $a_n/b_n \to 0$ as $n \to \infty$.

**Assumption 1** *The regressor $\phi_k$ is $\mathcal{F}_k$-measurable and bounded. Also, there is a known constant $L > 0$ such that $\|\theta\|_1 \leq L$.*

Under Assumption 1, there exists a positive bounded sequence $\{M_k\}$ such that $\max\limits_{1 \leq i \leq 2^{m_1}(3+3m_2)} |\phi_{k,i}| \leq M_k$ for all $k \geq 0$, where $\phi_{k,i}$ is the $i$-th component of $\phi_k$ and $M_k$ is bigger than the natural constant $e$. The time-varying bound $M_k$ will be used in the subsequent design of the algorithm.

**Assumption 2** *The thresholds $L_k$ and $N_k$ defined in (3) are $\mathcal{F}_k$-measurable, uniformly bounded with respect to the sampling path, and strictly ordered, satisfying $L_k < N_k$ for all $k$.*

Note that Assumption (2) is quite general for ensuring the saturated property of the function $S_k(\cdot)$.

**Assumption 3** *The random noise $\{w_k, \mathcal{F}_k\}$ is a martingale difference sequence, and satisfies*

$$\sup_{k \geq 0} E\big[|w_{k+1}|^4|\mathcal{F}_k\big] < \infty, \quad a.s. \tag{11}$$

*And the conditional expectation function $G_k(x) \triangleq E[S_k(x+w_{k+1})|\mathcal{F}_k]$ is differentiable, and its derivation function $G_k'(x)$ satisfies*

$$0 < \underline{g} = \inf_{|x|<C,k\geq0} G_k'(x) \leq \sup_{|x|<C,k\geq0} G_k'(x) = \overline{g} < \infty, \tag{12}$$

*where $C > 0$ is any constant.*

For convenience of analysis, we also introduce the following notations used in the subsequent sections:

$$\tilde{\theta}_k = \theta - \theta_k, \tag{13}$$

$$v_{k+1} = z_{k+1} - G_k(\phi_k^\tau \theta), \tag{14}$$

$$\psi_k = G_k(\phi_k^\tau \theta) - G_k(\phi_k^\tau \theta_k), \tag{15}$$

$$\underline{g}_k = \inf_{|x| \leq \max\{M_k L, M_k \|\theta_k\|_1\}} G_k'(x), \tag{16}$$

$$\overline{g}_k = \sup_{|x| \leq \max\{M_k L, M_k \|\theta_k\|_1\}} G_k'(x), \tag{17}$$

where $\theta_k, k \geq 0$ is the estimate for $\theta$, which is obtained by our algorithm introduced below.

### B. A Gradient-Based Two-Stage Algorithm

Now, we propose our two-stage gradient algorithm, the details of this algorithm are outlined in Algorithm 1 below.

In Algorithm 1, the Adam algorithm in the second stage is based on the following loss for the $h$-th batch of training data during the $l$-th epoch ($h = 1, \cdots, \lfloor n/\mathcal{T} \rfloor$, $l = 1, \cdots, N$):

$$\frac{1}{\mathcal{T}} \sum_{k=(h-1)\mathcal{T}+1}^{h\mathcal{T}} \frac{|z_k^{(l)} - \hat{z}_k^{(l)}|}{z_k^{(l)}} + \gamma \left| \frac{1}{\mathcal{T}} \sum_{k=(h-1)\mathcal{T}+1}^{h\mathcal{T}} \hat{e}_k^{(l)} - \bar{e} \right|, \tag{25}$$

where $z_k^{(l)}$ denotes the actual sentence of the $k$-th case in the $l$-th epoch, and $\hat{z}_k^{(l)} = S_k\Big( [a_k + \hat{b}_{h-1}^{(l)} x_k^{(1)} + \hat{c}_{h-1}^{(l)} x_k^{(2)}] \times \prod_{i=1}^{m_1} \big(1 + (\hat{p}_{h-1}^{(i)})^{(l)} v_k^{(i)}\big) \times \big[1 + \sum_{j=1}^{m_2} (\hat{q}_{h-1}^{(j)})^{(l)} u_k^{(j)} + \hat{e}_k^{(l)}\big]\Big)$ denotes the corresponding predictive sentence, $\hat{e}_k^{(l)} = \hat{\Gamma}_{h-1}^{(l)} \sigma\Big( \hat{B}_{h-1}^{(l)} \sigma\Big( \hat{A}_{h-1}^{(l)} \eta_k + (\hat{b}_{h-1}^{(1)})^{(l)} \Big) + (\hat{b}_{h-1}^{(2)})^{(l)} \Big) + (\hat{b}_{h-1}^{(3)})^{(l)}$ denotes the corresponding predictive time-varying bias term, where $\hat{b}_{h-1}^{(l)}$, $\hat{c}_{h-1}^{(l)}$, $(\hat{p}_{h-1}^{(i)})^{(l)}$, $(\hat{q}_{h-1}^{(j)})^{(l)}$, $\hat{\Gamma}_{h-1}^{(l)}$, $\hat{B}_{h-1}^{(l)}$, $\hat{A}_{h-1}^{(l)}$, $(\hat{b}_{h-1}^{(1)})^{(l)}$, $(\hat{b}_{h-1}^{(2)})^{(l)}$ and $(\hat{b}_{h-1}^{(3)})^{(l)}$ are the estimates of the true parameters in the model (1) after training on the $(h-1)$-th batch data in the $l$-th epoch. $\bar{e}$ is the ASG-based estimate of the

---

**Algorithm 1** Two-Stage Learning (TSL) Algorithm

---

1: **Input**: The temporally ordered training dataset $D_1 = \{\phi_k, z_{k+1}\}_{k=0}^{n-1}$ , the held-out testing dataset $D_2 = \{\phi_k, z_{k+1}\}_{k=1}^{n_2}$, the arbitrarily chosen initial estimates $\theta_0 \in \mathbb{R}^{p \times 1}$, $\hat{\Gamma}_0 \in \mathbb{R}^{1 \times m}$, $\hat{B}_0 \in \mathbb{R}^{m \times m}$, $\hat{A}_0 \in \mathbb{R}^{m \times m_3}$, $\hat{b}_0^{(1)} \in \mathbb{R}^m$, $\hat{b}_0^{(2)} \in \mathbb{R}^m$, $\hat{b}_0^{(3)} \in \mathbb{R}$, the epoch count of $N$, the batch size $\mathcal{T}$ ($\mathcal{T} \leq n$), hyper-parameters $\alpha > 1$, $\eta_1 > 0$, $\varepsilon > 0$, $\mu \in (0, 1]$ , $\beta_1 \in (0, 1)$ and $\beta_2 \in (0, 1)$.

2: **Output**: The final parameter estimate $\hat{\Theta}_{\lfloor n/\mathcal{T} \rfloor}^{(N)}$, the predictive sentences for sentencing cases in $D_2$.

3: **# Stage 1: ASG-based initialization for the unknown mechanistic parameter in the model (1)**.

4:     **for** $k = 0$ **to** $n - 1$ **do**

$$\theta_{k+1} = \theta_k + \frac{\mu \bar{g}_k \phi_k}{r_k^{\frac{1}{2}} \log^{\frac{\alpha}{2}} r_k} \left[ z_{k+1} - G_k(\theta_k^T \phi_k) \right], \tag{18}$$

$$r_k = M^4 p^2 + \sum_{i=1}^{k} \bar{g}_i^2 \|\phi_i\|^2, r_0 = M^4 p^2, \tag{19}$$

5:     **end for**

6:     Calculate the parameter estimates of the mechanistic part in (1) and the bias term by using the estimate $\theta_n$:

$$b_0 = \frac{\theta_n^{(2^{m_1}(1+m_2)+1)}}{\theta_n^{(1)}}, \quad c_0 = \frac{\theta_n^{(2^{m_1}(2+2m_2)+1)}}{\theta_n^{(1)}}, \quad \bar{e} = \theta_n^{(1)} - 1, \tag{20}$$

$$p_{i0} = \frac{\theta_n^{(i+1)}}{\theta_n^{(1)}}, \qquad q_{j0} = \theta_n^{(2^{m_1}j+1)}, \tag{21}$$

7: **# Stage 2: Adam-based estimation for all unknown parameters in the model (1)**.

8:     **for** $l = 1$ **to** $N$ **do**

9:         **if** $l = 1$ **then**

$$\hat{\Theta}_0^{(1)} = [b_0, c_0, p_{10}, \cdots, p_{m_10}, q_{10}, \cdots, q_{m_20}, \text{vec}(\hat{\Gamma}_0)^T, \text{vec}(\hat{B}_0)^T, \text{vec}(\hat{A}_0)^T, (\hat{b}_0^{(1)})^T, (\hat{b}_0^{(2)})^T, \hat{b}_0^{(3)}]^T,$$
$$m_0^{(1)} = 0, \ v_0^{(1)} = 0.$$

10:         **else**

$$\hat{\Theta}_0^{(l)} = \hat{\Theta}_{\lfloor n/\mathcal{T} \rfloor}^{(l-1)}, \ m_0^{(l)} = m_{\lfloor n/\mathcal{T} \rfloor}^{(l-1)}, \ v_0^{(l)} = v_{\lfloor n/\mathcal{T} \rfloor}^{(l-1)}.$$

11:         **end if**

12:         **for** $h = 1$ **to** $\lfloor n/\mathcal{T} \rfloor$ **do**

$$\hat{\Theta}_h^{(l)} = \hat{\Theta}_{h-1}^{(l)} + \frac{\eta_1}{\varepsilon + \sqrt{\frac{v_h^{(l)}}{1-\beta_2^h}}} \cdot \frac{m_h^{(l)}}{1 - \beta_1^h}, \tag{22}$$

$$m_h^{(l)} = \beta_1 m_{h-1}^{(l)} + (1 - \beta_1) g_h^{(l)}, \tag{23}$$

$$v_h^{(l)} = \beta_2 v_{h-1}^{(l)} + (1 - \beta_2)[g_h^{(l)}]^2, \tag{24}$$

13:         **end for**

14:     **end for**

15: Calculate the predictive sentences in $D_2$ by using the final estimate $\hat{\Theta}_{\lfloor n/\mathcal{T} \rfloor}^{(N)}$.

---

averaged bias term obtained by (20). The regularization coefficient $\gamma$ is determined to guide the neural network output towards matching the averaged bias term.

Moreover, as for other related notations in Algorithm 1, $p \triangleq 2^{m_1}(3 + 3m_2)$ in (19) is the dimension of unknown parameters $\theta$ defined in (5), $\theta_n^{(\xi)}$ ($\xi = 1, \cdots, p$) in (20)-(21) is the $\xi$-th component of the estimate $\theta_n$. $g_h^{(l)}$ in (23)-(24) denotes the negative gradient of the loss function (25) evaluated at $\hat{\Theta}_{h-1}^{(l)} = [\hat{b}_{h-1}^{(l)}, \hat{c}_{h-1}^{(l)}, (\hat{p}_{h-1}^{(1)})^{(l)}, \cdots, (\hat{p}_{h-1}^{(m_1)})^{(l)}, (\hat{q}_{h-1}^{(1)})^{(l)}, \cdots, (\hat{q}_{h-1}^{(m_2)})^{(l)}, \text{vec}(\hat{\Gamma}_{h-1}^{(l)})^T, \text{vec}(\hat{B}_{h-1}^{(l)})^T, \text{vec}(\hat{A}_{h-1}^{(l)})^T,$ $[(\hat{b}_{h-1}^{(1)})^{(l)}]^T, [(\hat{b}_{h-1}^{(2)})^{(l)}]^T, (\hat{b}_{h-1}^{(3)})^{(l)}]^T$ with $\text{vec}(\cdot)$ denoting the vectorization of a matrix. Besides, the addition, multiplication, and division operations in the Adam algorithm are performed element-wise.

The essence of Algorithm 1 is that the ASG algorithm in the first stage generates parameter estimates close to the prediction-error minimizer for the mechanism model, which are then used to initialize the Adam algorithm for the hybrid model in the second stage, thereby enhancing overall sentencing prediction performance. It is worth noting that, in order to ensure that the neural network serves only as a compensatory component for

the mechanistic model, the loss function (25) of the Adam algorithm incorporates a penalty on the deviation of the averaged neural network output from the bias term $\bar{e}$ estimated by the ASG algorithm in the first stage.

**Remark 4** *The ASG algorithm specified in (18)-(19) is an adaptive algorithm, i.e., the algorithm updates the parameter estimate $\theta_{k+1}$ using only the current online data $\{\phi_k, z_{k+1}\}$ and the current estimate $\theta_k$. The algorithm is motivated by the classical ASG algorithm studied in [20]- [22], but achieves a faster theoretical convergence rate of the averaged regret due to the use of a larger adaptation rate, which will be proven in the subsequent theoretical analysis.*

**Remark 5** *The estimate $\theta_k$ obtained by the ASG-type algorithm (18)-(19) is bounded. Moreover, according to Assumption 3 and the definition of $\underline{g}_k, \bar{g}_k$, it follows that*

$$\inf_{k \geq 0}\{\underline{g}_k\} > 0, \quad \sup_{k \geq 0}\{\bar{g}_k\} < \infty. \qquad (26)$$

**Remark 6** *For the proposed TSL algorithm, in Stage 1, the ASG method is proposed since it can provide good initial estimates for prediction in Stage 2. Specifically, the predictor based on the ASG method possesses a global asymptotic convergence property, as established in Section IV below, thereby providing good initial parameter estimates close to the prediction-error minimizer for the mechanism model—a guarantee typically absent in most machine learning algorithms. Moreover, the expanded regression vector in (4) is very high-dimensional in sentencing prediction, making stochastic-gradient–type methods particularly suitable for handling such cases, since the computational load of the gradient algorithms is much lower than that of the Newton-type algorithms. In Stage 2, once the bias term $e$ is replaced by a neural network, the Adam algorithm is applied to improve the overall predictive capability of the model, as Adam is widely recognized as one of the most appropriate optimization algorithms for neural networks [33]. The advantages of using Adam are also confirmed by the experimental results presented in Section V.*

## IV. PREDICTION THEORY OF THE ASG ALGORITHM

In this section, we establish the global asymptotic convergence of the ASG-based predictor without requiring any excitation data conditions.

To facilitate the theoretical analysis, we first introduce the corresponding best predictor. By (10) and the definition of $G_k(\cdot)$, one can deduce that the best prediction of $z_{k+1}$ given $\mathscr{F}_k$ in the mean square sense is as follows:

$$E[z_{k+1}|\mathscr{F}_k] = G_k(\theta^T \phi_k). \qquad (27)$$

Since $\theta$ is unknown *a priori*, we replace the unknown parameter $\theta$ by its estimates $\theta_k$ and define the adaptive prediction for $z_{k+1}$ at time $k$ as follows:

$$\hat{z}_{k+1} = G_k(\theta_k^T \phi_k). \qquad (28)$$

From the above, we define the difference between the best prediction and the adaptive prediction for the saturated sentence $z_{k+1}$ as "regret", which can be expressed as follows:

$$R_k = [E[z_{k+1}|\mathscr{F}_k] - \hat{z}_{k+1}]^2 = [\psi_k]^2. \qquad (29)$$

Naturally, we expect the regret for $z_{k+1}$ to decrease as $k$ increases, ideally vanishing to zero. The following theorem shows that this can be achieved in the averaged sense without requiring any data excitation condition, such as the i.i.d. condition, etc.

**Theorem 1** *Under Assumptions 1-3, the accumulated regrets have the following upper bounds:*

$$\sum_{k=0}^{n-1} R_k = o\left(n^{\frac{1}{2}} \log^{\frac{\alpha}{2}} n\right), \quad a.s. \qquad (30)$$

where $r_k$ and $\alpha$ are defined in Algorithm 1.

**Corollary 1** *If the bounded condition on $\phi_{k,i}$ in Assumption 1 is relaxed to $|\phi_{k,i}| \leq Mk^\epsilon, 0 < \epsilon < \frac{1}{2}$ for any $1 \leq i \leq p$, and $\bar{g}_k$ and $r_k$ in the algorithm (18)-(19) are replaced by $\bar{g}_k = \sup\limits_{|x| \leq \max\{LMk^\epsilon, Mk^\epsilon\|\theta_k\|_1\}} G'_k(x)$ and $r_k = M^4 k^{4\epsilon} p^2 + \sum_{i=1}^{k} \bar{g}_i^2 \|\phi_i\|^2$, respectively, then the accumulated regrets possess the following upper bound:*

$$\sum_{k=0}^{n-1} R_k = o\left(n^{\frac{1}{2}+\epsilon} \log^{\frac{\alpha}{2}} n\right), \quad a.s. \qquad (31)$$

It can been seen that the averaged regrets $\frac{1}{n}\sum_{k=0}^{n-1} R_k$ in both Theorem 1 and Corollary 1 will converge to zero almost surely as $n \to \infty$. The proofs of the above results will be shown in the Appendix below.

## V. SENTENCING EXPERIMENTS

In this section, we demonstrate the superiority of the proposed model and algorithm based on a real-world CII dataset. First, the experimental settings and the baseline models for comparison are described, followed by a detailed analysis of the results.

### A. Experimental Settings

We conduct judicial sentencing experiments based on an available CII real-world dataset obtained from China Judgements Online[3], which contains 87,588 original minor injury judgment documents, as well as 9,228 original serious injury judgment documents from the

---

[3]https://wenshu.court.gov.cn/

period 2019 to 2024. These data are processed through feature extraction and quantification to facilitate subsequent experiments.

We provide a detailed explanation of the output boundary selection and feature interpretation for the SMNN model. Regarding the output boundaries, Article 234 of Chinese Criminal Law prescribes a fixed-term imprisonment ranging from six months to three years for minor injury cases and three years to ten years for serious injury cases, so we let the statutory penalty ranges $L_k \equiv 6$, $N_k \equiv 36$ and $L_k \equiv 36$, $N_k \equiv 120$ for minor and serious cases in (3), respectively (Unit: Month). According to the sentencing guidelines for the CII, penalty-determining factors $x_k^{(i)}, i = 1, 2$ in the model (1) represent the number of seriously and minorly injured victims, respectively. Primary sentencing adjust factors $v_k^{(i)}, i = 1, \cdots, 13$ in the model (1) include "juveniles aged 16–18 years," "juveniles aged 12–16 years," "elderly individuals over 75 years of age," "mentally disordered persons with diminished criminal responsibility," "deaf-mute and blind individuals," "excessive self-defense," "excessive act of necessity," "preparatory acts for crime," "attempted crime," "voluntary cessation of crime," "accessories," "coerced accomplices," and "instigators." Other sentencing adjust factors $u_k^{(j)}, j = 1, \cdots, 22$ in the model (1) include "grade I serious injury cases," "grade II serious injury cases," "grade I minor injury cases," "grade II minor injury cases," "voluntary surrender," "confession," "criminal reconciliation," "active compensation," "victim pardon," "principal offender," "crimes targeting vulnerable groups," "victim's contributory fault," "civil dispute-related offenses," "recidivism," "prior criminal records," "first-time offenders," "courtroom confession," "plea agreement acceptance," "armed affray," "mutual combat," "meritorious service," and "probation," among others.

The prediction accuracy metric adopts a relative accuracy with discretion (RAD), which is defined as follows:

$$\text{RAD} = 1 - \frac{1}{n_2} \sum_{k=1}^{n_2} \frac{\tilde{z}_k}{z_k} I(\tilde{z}_k > \max\{20\% z_k, 2\}), \quad (32)$$

where $\tilde{z}_k = |z_k - \hat{z}_k|$ with $z_k$ denoting the actual sentence for the $k$-th judicial case, and $\hat{z}_k$ representing its predicted value. $n_2$ is the total number of criminal cases in the testing set. The threshold $\max\{20\% z_k, 2\}$ represents the degree of adjustment allowed during the judges' sentencing process, indicating the judicial discretion. Here the component $20\% z_k$ is based on the sentencing guidelines [4], while the fixed value 2 is derived from interviews with judges and reflects practical discretion in real-world judicial decision-making. Compared to the prediction accuracy metrics employed in previous related

---

[4]See https://www.court.gov.cn/.

studies (e.g., classification metrics in [1]- [2], [6]), the metric (32) is both legally compatible and practically reasonable.

Moreover, the sentencing starting point $a_k$ in the model (1) is set as the lower bound of the sentencing starting range prescribed in law based on a series of comparative experiments(see [11]). The dataset is split into a training set and a testing set with a ratio of $4 : 1$. Some hyper-parameters for the ASG algorithm are configured as follows: the constant $\alpha$ in Algorithm 1 is set to 1.02, $\mu$ and $\bar{g}_k$ are set to 1, and the random noise sequence $\{w_k\}$ is assumed to be i.i.d. following a normal distribution $N(0, 25)$.

Additionally, other common hyper-parameters for the Adam algorithm involved in relevant experiments are set as follows: a batch size of 245, an epoch count of 30, a learning rate of 0.001, exponential decay rates for the first and second moment estimates of $\beta_1 = 0.9$, $\beta_2 = 0.999$, a smoothing coefficient of $\varepsilon = 10^{-8}$. The hidden layer is configured with 128 nodes. Importantly, during training across different epochs, data are fed in chronological order without shuffling to preserve the temporal structure inherent in the dataset. The regularization coefficient $\gamma$ in (25) is set to 0.2 and 1.4 for minor and serious injury cases respectively, based on comparative experiments of prediction accuracy. All experiments involving random initialization are conducted 10 times, with the best-performing initialization selected to optimize overall performance.

### B. Baseline Models

*1) SM Model with ASG-Based Algorithm:* First, we consider the SM model (see [31]), where the bias term in (1) is a fixed unknown constant $e$. Specifically, the model is as follows.

$$\begin{aligned} z_{k+1} = S_k\bigg( & \left[a_k + bx_k^{(1)} + cx_k^{(2)}\right] \\ & \times \prod_{i=1}^{m_1}(1 + p_i v_k^{(i)}) \\ & \times \left[1 + \sum_{j=1}^{m_2} q_j u_k^{(j)} + e\right] + w_{k+1}\bigg). \end{aligned} \quad (33)$$

The proposed ASG algorithm (18)–(19) will be used for optimization based on the model (33).

*2) SNN Model with Adam-Based Algorithm:* Second, we conduct experiments utilizing a saturated neural network (SNN) model (see [12]) specified as follows.

$$\begin{aligned} & z_k \\ & = S_k\left[W^{(3)}\sigma\big(W^{(2)}\sigma(W^{(1)}X_k + B^{(1)}) + B^{(2)}\big) + B^{(3)}\right], \end{aligned} \quad (34)$$

where $W^{(i)}, B^{(i)}, i = 1, 2, 3$ are unknown weighting parameters, and $X_k$ is the input data composed of various sentencing factors. Other notations are consistent with those defined in Section II. The Adam algorithm (22)-(24) will be used for optimization based on the model (34).

*3) SMNN Model with Adam-Based Algorithm:* Third, we train our SMNN model (1) by Adam optimizer (22)–(24) with all the parameters being initialized randomly.

*4) SMNN Model with TSL Algorithm:* Fourth, we train the SMNN model (1) by our TSL optimizer described in Algorithm 1, where the mechanistic parameters are initialized using the ASG algorithm in the first stage, and other parameters are initialized randomly.

### C. Comparison Experiments on Sentencing Prediction

To further improve the sentencing prediction accuracy for CII in the existing literature (e.g., [11], [31], [32]), we perform comparison experiments to demonstrate the improved accuracy of our SMNN model and the TSL algorithm compared to the above baseline methods using the same available data.

Figs. 1-2 illustrate the prediction accuracy trends of our method in comparison to other related known methods on the testing set for minor and serious cases, respectively. It can be shown that:

- The SMNN model with the ASG-Adam algorithm consistently achieves the highest prediction accuracy of 86.66% and 95.25% for minor and serious cases respectively, surpassing the SM model with the ASG algorithm (80.69%; 90.59%) and the SNN model with the Adam algorithm (85.66%; 94.38%). This demonstrates the advantages of combining both the SM model with the NN model as well as integrating the ASG algorithm with the Adam algorithm, allowing for more accurate sentencing prediction.
- The SMNN model with our ASG-Adam algorithm outperforms the same model with the Adam algorithm with random initialization (86.46%; 93.29%), which demonstrates that our ASG algorithm (18)-(19) in the first stage provides good initial parameter estimates for the Adam optimizer (22)-(24) in the second stage, thereby improving prediction accuracy. This advantage stems from the ASG-based initialization, which guides the optimization process toward the global error minimum. In contrast, random initialization requires exploring a broader solution space and is more prone to local error minima, often resulting in inferior prediction accuracy.
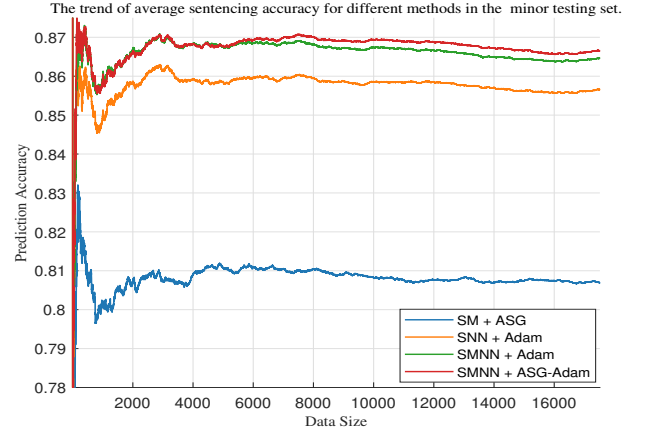


Fig. 1. The trend of average sentencing accuracy of minor cases for different methods in the testing set.
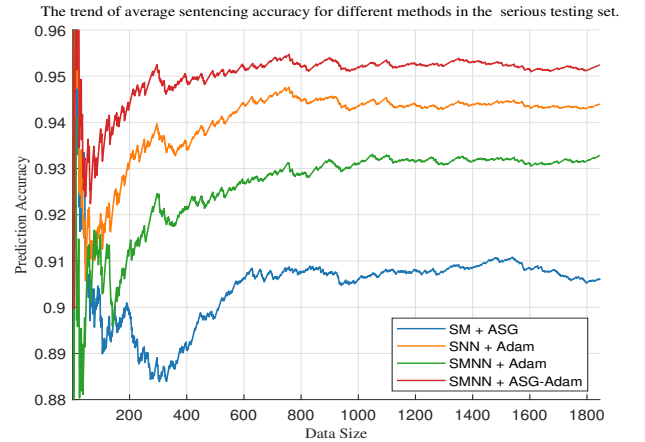


Fig. 2. The trend of average sentencing accuracy of serious cases for different methods in the testing set.

## VI. CONCLUSION

Motivated by the imperative demand in judicial practice for highly accurate and reliable sentencing prediction, this paper proposed a SMNN hybrid model, which integrates the sentencing logic model with a neural network. The ASG algorithm is applied to update the parameters of the mechanistic component within the SMNN model, after which the Adam algorithm is utilized to further optimize the SMNN model with the aim of enhancing predictive accuracy. Both theoretical analysis and sentencing experiments demonstrate that the estimates produced by the ASG algorithm in the first stage can provide good enough initial values for the implementation of the prediction algorithm in the second stage. Moreover, sentencing experiments also reveal that the neural networks introduced in the hybrid model is indeed helpful for improving the prediction accuracy in

the second stage. For future investigation, it would be interesting to establish the global convergence of the Adam algorithm in the second stage, and to apply our prediction algorithms to other judicial crimes other than the intentional injury studied here.

## VII. APPENDIX

**Lemma 1** ( [23] Theorem 1.3.2) *Let $\{f_k, \mathscr{F}_k\}$ and $\{\alpha_k, \mathscr{F}_k\}$ be two non-negative adapted sequences. If $E[f_{k+1}|\mathscr{F}_k] \leq f_k + \alpha_k$, a.s. for any $k \geq 1$ and $\sum_{i=1}^{\infty} \alpha_i < \infty$, a.s., then $f_k$ will converge to a finite limit a.s.*

**Lemma 2** ( [23] Theorem 1.2.15) *Let $D_k = C + \sum_{j=1}^{k} d_j, d_j \geq 0$ and $D_0 = C$ with $C > 1$ being any constant, then we have*

$$\sum_{j=1}^{\infty} \frac{d_j}{D_j \log^{\alpha} D_j} < \infty, \forall \alpha > 1. \tag{35}$$

**Lemma 3** ( [25]) *Let $\{\omega_n, \mathcal{F}_n\}$ be a martingale difference sequence and $\{f_n, \mathcal{F}_n\}$ an adapted sequence. If*

$$\sup_n \mathbb{E}\left[|\omega_{n+1}|^{\alpha} \mid \mathcal{F}_n\right] < \infty, \quad a.s. \tag{36}$$

*for some $\alpha \in (0, 2]$, then as $n \to \infty$, we have $\forall \eta > 0$,*

$$\sum_{i=0}^{n} f_i \omega_{i+1} = O\left(s_n(\alpha) \log^{\frac{1}{\alpha}+\eta}\left(s_n^{\alpha}(\alpha) + e\right)\right) \quad a.s., \tag{37}$$

*where $s_n(\alpha) = \left(\sum_{i=0}^{n} |f_i|^{\alpha}\right)^{\frac{1}{\alpha}}$.*

Inspired by the analysis ideas mentioned in [23], [17], and [18], etc, the proof of Theorem 1 is as follows:

**Proof of Theorem 1.** By the algorithm (18)-(19) and the differential mean value theorem, we have

$$\tilde{\theta}_{k+1} = \tilde{\theta}_k - \frac{\mu \bar{g}_k \phi_k}{r_k^{\frac{1}{2}} \log^{\frac{\alpha}{2}} r_k}\left[G_k'(\xi_k)\phi_k^T \tilde{\theta}_k + v_{k+1}\right]$$

$$= \left(I - \frac{\mu \bar{g}_k G_k'(\xi_k)\phi_k \phi_k^T}{r_k^{\frac{1}{2}} \log^{\frac{\alpha}{2}} r_k}\right)\tilde{\theta}_k - \frac{\mu \bar{g}_k \phi_k}{r_k^{\frac{1}{2}} \log^{\frac{\alpha}{2}} r_k}v_{k+1}, \tag{38}$$

where $\xi_k \in (\min\{\theta^\tau \phi_k, \theta_k^\tau \phi_k\}, \max\{\theta^\tau \phi_k, \theta_k^\tau \phi_k\})$. By this and the elementary inequality we have

$$\tilde{\theta}_{k+1}^T \tilde{\theta}_{k+1}$$

$$= \left\{\left(I - \frac{\mu \bar{g}_k G_k'(\xi_k)\phi_k \phi_k^T}{r_k^{\frac{1}{2}} \log^{\frac{\alpha}{2}} r_k}\right)\tilde{\theta}_k - \frac{\mu \bar{g}_k \phi_k}{r_k^{\frac{1}{2}} \log^{\frac{\alpha}{2}} r_k}v_{k+1}\right\}^T$$

$$\left\{\left(I - \frac{\mu \bar{g}_k G_k'(\xi_k)\phi_k \phi_k^T}{r_k^{\frac{1}{2}} \log^{\frac{\alpha}{2}} r_k}\right)\tilde{\theta}_k - \frac{\mu \bar{g}_k \phi_k}{r_k^{\frac{1}{2}} \log^{\frac{\alpha}{2}} r_k}v_{k+1}\right\}$$

$$= \tilde{\theta}_k^T \tilde{\theta}_k - 2\frac{\mu \bar{g}_k G_k'(\xi_k)\left[\tilde{\theta}_k^T \phi_k\right]^2}{r_k^{\frac{1}{2}} \log^{\frac{\alpha}{2}} r_k} - 2\frac{\mu \bar{g}_k \tilde{\theta}_k^T \phi_k}{r_k^{\frac{1}{2}} \log^{\frac{\alpha}{2}} r_k}v_{k+1}$$

$$+ \frac{\mu^2 \bar{g}_k^2 [G_k'(\xi_k)]^2 \|\phi_k\|^2 \left[\tilde{\theta}_k^T \phi_k\right]^2}{r_k \log^{\alpha} r_k}$$

$$+ 2\frac{\mu^2 \bar{g}_k^2 G_k'(\xi_k)\|\phi_k\|^2 \tilde{\theta}_k^T \phi_k}{r_k \log^{\alpha} r_k}v_{k+1} + \frac{\mu^2 \bar{g}_k^2 \|\phi_k\|^2}{r_k \log^{\alpha} r_k}v_{k+1}^2$$

$$\leq \tilde{\theta}_k^T \tilde{\theta}_k - 2\frac{\mu \bar{g}_k G_k'(\xi_k)\left[\tilde{\theta}_k^T \phi_k\right]^2}{r_k^{\frac{1}{2}} \log^{\frac{\alpha}{2}} r_k} - 2\frac{\mu \bar{g}_k \tilde{\theta}_k^T \phi_k}{r_k^{\frac{1}{2}} \log^{\frac{\alpha}{2}} r_k}v_{k+1}$$

$$+ (1 + b)\frac{\mu^2 \bar{g}_k^2 [G_k'(\xi_k)]^2 \|\phi_k\|^2 \left[\tilde{\theta}_k^T \phi_k\right]^2}{r_k \log^{\alpha} r_k}$$

$$+ \left(1 + \frac{1}{b}\right)\frac{\mu^2 \bar{g}_k^2 \|\phi_k\|^2}{r_k \log^{\alpha} r_k}v_{k+1}^2, \tag{39}$$

where $0 < b < 1$ is a constant. since $S_k'(x) \leq 1, \forall x \in \mathbb{R}$ implies $0 \leq \bar{g}_k G_k'(\xi_k) \leq 1$, summing up both sides of (39) from $k = 0$ to $n - 1$, and noticing $\|\phi_k\|^2 < r_k^{\frac{1}{2}} \log^{\frac{\alpha}{2}} r_k$ for all $k \geq 0$, we know that

$$\|\tilde{\theta}_n\|^2 + (1 - b)\sum_{k=0}^{n-1} \frac{\mu \bar{g}_k G_k'(\xi_k)\left[\tilde{\theta}_k^T \phi_k\right]^2}{r_k^{\frac{1}{2}} \log^{\frac{\alpha}{2}} r_k}$$

$$= O\left(\sum_{k=0}^{n-1} \frac{\bar{g}_k \tilde{\theta}_k^T \phi_k}{r_k^{\frac{1}{2}} \log^{\frac{\alpha}{2}} r_k}v_{k+1}\right) + O\left(\sum_{k=0}^{n-1} \frac{\bar{g}_k^2 \|\phi_k\|^2}{r_k \log^{\alpha} r_k}v_{k+1}^2\right). \tag{40}$$

Now we analyze the RHS of (40) term by term. Note that by (11), the elementary inequality $2ab \leq a^2 + b^2$ and the inequality $E^4[|x||\mathscr{F}_k] \leq E[|x|^4|\mathscr{F}_k]$, we have

$$E\left[|v_{k+1}|^4 \mid \mathscr{F}_k\right]$$

$$= O\left(E\left[\left|S_k(\theta^\tau \phi_k + w_{k+1}) - S_k(\theta^\tau \phi_k)\right|^4 \big| \mathscr{F}_k\right]\right)$$

$$+ O\left(E\left[\left|S_k(\theta^\tau \phi_k) - E[S_k(\theta^\tau \phi_k + w_{k+1})|\mathscr{F}_k]\right|^4 \big| \mathscr{F}_k\right]\right)$$

$$= O\left(E\left[\left|S_k(\theta^\tau \phi_k + w_{k+1}) - S_k(\theta^\tau \phi_k)\right|^4 \big| \mathscr{F}_k\right]\right) + O(1)$$

$$= O\left(E\left[|w_{k+1}|^4|\mathscr{F}_k\right]\right) + O(1) = O(1), \quad a.s. \tag{41}$$

Note that by the definition of $v_{k+1}$, we have

$$E\left(v_{k+1} \mid \mathscr{F}_k\right)$$

$$= E\left[S_k\left(\theta^\tau \phi_k + w_{k+1}\right) \mid \mathscr{F}_k\right]$$

$$- E\left[E\left[S_k\left(\theta^\tau \phi_k + w_{k+1}\right) \mid \mathscr{F}_k\right] \mid \mathscr{F}_k\right] = 0. \tag{42}$$

So by (41), (42) and Lemma 3, we have

$$\sum_{k=0}^{n-1} \frac{\bar{g}_k \tilde{\theta}_k^T \phi_k}{r_k^{\frac{1}{2}} \log^{\frac{\alpha}{2}} r_k}v_{k+1} = o\left(\sum_{k=0}^{n-1} \frac{\bar{g}_k^2 \left[\tilde{\theta}_k^T \phi_k\right]^2}{r_k^{\frac{1}{2}} \log^{\frac{\alpha}{2}} r_k}\right) + O(1), a.s. \tag{43}$$

Moreover, by (41), (42) and Lemma 2 and Lemma 3, we have

$$\sum_{k=0}^{n-1} \frac{\bar{g}_k^2 \|\phi_k\|^2}{r_k \log^\alpha r_k} v_{k+1}^2$$

$$\leq \sum_{k=0}^{n-1} \frac{\bar{g}_k^2 \|\phi_k\|^2}{r_k \log^\alpha r_k} \left[ v_{k+1}^2 - E[v_{k+1}^2 | \mathscr{F}_k] \right]$$

$$+ \sup_{k \geq 0} E[v_{k+1}^2 | \mathscr{F}_k] \cdot \sum_{k=0}^{n-1} \frac{\bar{g}_k^2 \|\phi_k\|^2}{r_k \log^\alpha r_k}$$

$$= O\left( \sum_{k=0}^{n-1} \frac{\bar{g}_k^2 \|\phi_k\|^2}{r_k \log^\alpha r_k} \right) + O(1)$$

$$= O(1), \quad a.s. \tag{44}$$

Combining (43)-(44) with (40), we have

$$\sum_{k=0}^{n-1} \frac{\bar{g}_k G_k'(\xi_k) \left[ \tilde{\theta}_k^T \phi_k \right]^2}{r_k^{\frac{1}{2}} \log^{\frac{\alpha}{2}} r_k} = O(1), \quad a.s. \tag{45}$$

Note that since $\phi_k$ is bounded, we have $r_k = O(k)$. Moreover, if $r_k \to \infty$ as $k \to \infty$, then we have by (45) and the Kronecker Lemma, we obtain

$$\sum_{k=0}^{n-1} R_k = o\left( n^{\frac{1}{2}} \log^{\frac{\alpha}{2}} n \right), \quad a.s. \tag{46}$$

Otherwise, if $r_k \not\to \infty$, then denominator of (45) is bounded, and (46) is also satisfied.

**Proof of Remark 5.** By (18), we have

$$\tilde{\theta}_{k+1}^T \tilde{\theta}_{k+1} = \left[ \tilde{\theta}_k - \frac{\mu \bar{g}_k \phi_k}{r_k^{\frac{1}{2}} \log^{\frac{\alpha}{2}} r_k} \left[ z_{k+1} - G_k \left( \phi_k^T \theta_k \right) \right] \right]^T$$

$$\left[ \tilde{\theta}_k - \frac{\mu \bar{g}_k \phi_k}{r_k^{\frac{1}{2}} \log^{\frac{\alpha}{2}} r_k} \left[ z_{k+1} - G_k \left( \phi_k^T \theta_k \right) \right] \right]. \tag{47}$$

Set $\sup_{k \geq 0} \max\{|L_k|, |N_k|\} \leq U < \infty$, from (14) and (15) we know that

$$\left\| \tilde{\theta}_{k+1} \right\|^2$$

$$= \left\| \tilde{\theta}_k \right\|^2 - 2 \frac{\mu \bar{g}_k}{r_k^{\frac{1}{2}} \log^{\frac{\alpha}{2}} r_k} \left[ z_{k+1} - G_k \left( \phi_k^T \theta_k \right) \right] \tilde{\theta}_k^T \phi_k$$

$$+ \frac{\mu^2 \bar{g}_k^2 \left[ z_{k+1} - G_k \left( \phi_k^T \theta_k \right) \right]^2}{r_k \log^\alpha r_k} \|\phi_k\|^2$$

$$\leq \left\| \tilde{\theta}_k \right\|^2 - 2 \frac{\mu \bar{g}_k}{r_k^{\frac{1}{2}} \log^{\frac{\alpha}{2}} r_k} \left[ v_{k+1} + \psi_k \right] \tilde{\theta}_k^T \phi_k$$

$$+ \frac{4U^2 \bar{g}_k^2}{r_k \log^\alpha r_k} \|\phi_k\|^2. \tag{48}$$

Take the conditional expectation for both sides of (48), and by (42) and differential mean value theorem, we know that there exists a random variable $\xi_k \in (\min\{\theta^\tau \phi_k, \theta_k^\tau \phi_k\}, \max\{\theta^\tau \phi_k, \theta_k^\tau \phi_k\})$ such that

$$E\left[ \left\| \tilde{\theta}_{k+1} \right\|^2 | \mathscr{F}_k \right]$$

$$\leq \left\| \tilde{\theta}_k \right\|^2 - 2 \frac{\mu \bar{g}_k G_k'(\xi_k)(\tilde{\theta}_k^T \phi_k)^2}{r_k^{\frac{1}{2}} \log^{\frac{\alpha}{2}} r_k} + \frac{4U^2 \bar{g}_k^2 \|\phi_k\|^2}{r_k \log^\alpha r_k}$$

$$\leq \left\| \tilde{\theta}_k \right\|^2 + \frac{4U^2 \bar{g}_k^2 \|\phi_k\|^2}{r_k \log^\alpha r_k}. \tag{49}$$

Moreover, by Lemma 2 we obtain $\sum_{k=1}^{n} \frac{\bar{g}_k^2 \|\phi_k\|^2}{r_k \log^\alpha r_k} = O(1)$, by this, (49) and Lemma 1, we know that there exists a constant $S \geq 0$ such that $\|\tilde{\theta}_k\| \to S < \infty$ as $k \to \infty$. Combining this with the boundedness of $\theta$, we can easily obtain the estimate $\theta_k$ is bounded.

**Proof of Corollary 1.** Note that $\|\phi_k\|^2 \leq r_k^{\frac{1}{2}} \log^{\frac{\alpha}{2}} r_k$, similar to the analysis of (40) and (43)-(44), we have

$$\sum_{k=0}^{n-1} \bar{g}_k G_k'(\xi_k) \left[ \tilde{\theta}_k^T \phi_k \right]^2 = o(n^{\frac{1}{2}+\epsilon} \log^{\frac{\alpha}{2}} n), a.s., \tag{50}$$

which implies that (31) holds.

## REFERENCES

[1] Kim, Y. "Convolutional neural networks for sentence classification," *Proceedings of Conference on Empirical Methods in Natural Language Processing*, 2014, pp. 1746–1751.

[2] Zhong, H., Guo, Z., Tu, C., et al. "Legal judgment prediction via topological learning," *Proceedings of Conference on Empirical Methods in Natural Language Processing*, 2018, pp. 3540–3549.

[3] Yang, W., Jia, W., Zhou, X., et al. "Legal judgment prediction via multi-perspective bi-feedback network," *Proceedings of IEEE Joint International Information Technology and Artificial Intelligence Conference*, 2019, pp. 4085–4091.

[4] Dong, Q., and Niu, S. "Legal judgment prediction via relational learning," *Proceedings of International ACM SIGIR Conference on Research and Development in Information Retrieval*, 2021, pp. 983–992.

[5] Xu, N., Wang, P., Chen, L., et al. "Distinguish confusing law articles for legal judgment prediction," *Proceedings of the Annual Meeting of the Association for Computational Linguistics*, 2020, pp. 3086–3095.

[6] Yue, L., Liu, Q., Jin, B., et al. "Neurjudge: A circumstance-aware neural framework for legal judgment prediction," *Proceedings of the International ACM SIGIR Conference on Research and Development in Information Retrieval*, 2021, pp. 973–982.

[7] Zhong, H., Wang, Y., Tu, C., et al. "Iteratively questioning and answering for interpretable legal judgment prediction," *Proceedings of the AAAI Conference on Artificial Intelligence*, 2020, pp. 1250–1257.

[8] Yang, S., Tong, S., Zhu, G., et al. "MVE-FLK: A multi-task legal judgment prediction via multi-view encoder fusing legal keywords," *Knowledge-Based Systems*, 2022, **239**: 107960.

[9] Zhao, Q., Gao, T., and Guo, N. "LA-MGFM: A legal judgment prediction method via semi-enhanced graph neural networks and multi-graph fusion mechanism," *Information Processing & Management*, 2023, **60**(5): 103455.

[10] Li, L., Liu, D., Zhao, L., et al. "Evidence mining for interpretable charge prediction via prompt learning," *IEEE Transactions on Computational Social Systems*, 2024, **11**(4): 4556–4566.

[11] Wang, F., Zhang, L., & Guo, L. "Applications of nonlinear recursive identification theory in the analyses of sentencing data," *Scientia Sinica Informationis*, **52**(10), 1837–1852, 2022.

[12] Dai, R., Wang, F., & Guo, L. "A New Adaptive Prediction Algorithm for Judicial Sentencing with Empirical Studies," *Journal of Systems Science and Complexity*, **38**(1), 3-20, 2025.

[13] Devlin, J., Chang, M., Lee, K., et al. "BERT: Pre-training of deep bidirectional transformers for language understanding," *Proceedings of Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 2019, pp. 4171–4186.

[14] Loshchilov, I., Hutter, F. "Decoupled Weight Decay Regularization," *Proceedings of International Conference on Learning Representations*, 2019.

[15] Kingma, D. P., and Ba, J. "Adam: A method for stochastic optimization," *Proceedings of International Conference on Learning Representations (ICLR)*, 2015.

[16] Zeiler, M. D. "Adadelta: an adaptive learning rate method," *arXiv preprint arxiv:1212.5701*, 2012.

[17] Zhang, L., & Guo, L. "Adaptive identification with guaranteed performance under saturated observation and non-persistent excitation," *IEEE Transactions on Automatic Control*, **69**(3), 1584–1599, 2023.

[18] Zheng, X., & Guo, L. "$L_1$-based adaptive identification with saturated observations," *IEEE Transactions on Automatic Control*, **70**(9), 5836-5847, 2025.

[19] Zhang, L., & Guo, L. "Adaptive tracking control with binary-valued output observations," *arXiv preprint*, arXiv:2411.05975, 2024.

[20] Chen, H. F., & Guo, L. "Consistency of parameter estimates for discrete-time linear systems," *Journal of Systems Science and Mathematical Sciences*, **5**(2), 81–93, 1985.

[21] Chen, H. F., & Guo, L. "Strong consistency of recursive identification without persistent excitation condition," *Acta Mathematicae Applicatae Sinica*, **2**(2), 133–145, 1985.

[22] Chen, H. F., & Guo, L. "The limit of stochastic gradient algorithm for identifying systems excited non-persistently," *Kexue Tongbao (Science Bulletin)*, 6–9, 1986.

[23] Guo, L. *Time-varying stochastic systems: Stability and adaptive theory* (2nd ed.), Science Press, Beijing, China, 2020.

[24] Chen, H. F. *Stochastic approximation and its applications*, Kluwer Academic Publishers, Dordrecht, Netherlands, 2002.

[25] Chen, H. F., & Guo, L. *Identification and stochastic adaptive control*, Birkhäuser, Boston, MA, 1991.

[26] Sun, J., Kim, Y. W., & Wang, L. "Aftertreatment control and adaptation for automotive lean burn engines with HEGO sensors," *International Journal of Adaptive Control and Signal Processing*, **18**(2), 145–166, 2008.

[27] Appadwedula, S., Veeravalli, V. V., & Jones, D. L. "Decentralized detection with censoring sensors," *IEEE Transactions on Signal Processing*, **56**(4), 1362–1373, 2008.

[28] Tobin, J. "Estimation of relationships for limited dependent variables," *Econometrica: Journal of the Econometric Society*, **26**(1), 24–36, 1958.

[29] Jeon, M. S., & Lee, J. H. "Estimation of willingness-to-pay for premium economy class by type of service," *Journal of Air Transport Management*, **84**, 28–35, 2020.

[30] Clark, T. G., Bradburn, M. J., Love, S. B., & Altman, D. G. "Survival analysis part I: Basic concepts and first analyses," *British Journal of Cancer*, **89**(2), 232–238, 2003.

[31] Guo, L. "The Integration of Law and Cybernetics in the Digital Age: An Exploration of Sentencing Research," *The 3rd Annual Conference on Computational Law*, Weihai, China, October 13, 2024.

[32] Jin, Y., Zheng, X., & Guo, L. "Adaptive sentencing prediction with guaranteed accuracy and legal interpretability," *arXiv preprint arXiv:2505.14011*, 2025. Available: https://doi.org/10.48550/arXiv.2505.14011

[33] Schmidt, R. M., Schneider, F., & Hennig, P. "Descending through a crowded valley—Benchmarking deep learning optimizers," *Proceedings of the International Conference on Machine Learning*, PMLR, 2021, pp. 9367-9376.