

Table 3: Assumption 5.1 verification.

	$m$	100	400	900	1600	2500
Mushroom	$g(m)$	$4.55 \pm 0.35$	$2.03 \pm 0.15$	$1.16 \pm 0.06$	$0.85 \pm 0.05$	$0.65 \pm 0.04$
Statlog	$g(m)$	$6.93 \pm 0.75$	$3.16 \pm 0.26$	$2.05 \pm 0.12$	$1.48 \pm 0.09$	$1.13 \pm 0.05$

## A Details of Theoretical Analysis

For theoretical analysis, we define the  $L$ -layered fully-connected neural network  $f$  as follows:

$$(A.1) \quad \begin{aligned} \mathbf{f}_1 &= \mathbf{W}_1 \mathbf{x}, \\ \mathbf{f}_l &= \mathbf{W}_l \text{ReLU}(\mathbf{f}_{l-1}), \quad 2 \leq l \leq L-1, \\ f(\mathbf{x}, \mathbf{w}) &= \sqrt{m} \mathbf{W}_L \text{ReLU}(\mathbf{f}_{L-1}), \end{aligned}$$

where  $\text{ReLU}(\cdot) := \max\{\cdot, 0\}$ ,  $\mathbf{x} \in \mathbb{R}^d$ ,  $\mathbf{W}_1 \in \mathbb{R}^{m \times d}$ ,  $\mathbf{W}_l \in \mathbb{R}^{m \times m}$ ,  $2 \leq l \leq L$ ,  $\mathbf{W}_L \in \mathbb{R}^{1 \times m}$ , and  $\mathbf{w} = (\text{vec}(\mathbf{W}_1); \dots; \text{vec}(\mathbf{W}_L)) \in \mathbb{R}^p$  is the collection vector of all the network parameters,  $p = dm + m^2(L-2) + m$ . The parameters are initialized as  $\mathbf{W}_l = (\mathbf{W}, \mathbf{0}; \mathbf{0}, \mathbf{W})$  for each  $1 \leq l \leq L-1$ , where each entry of  $\mathbf{W}$  is generated independently from  $\mathcal{N}(0, 4/m)$ , and  $\mathbf{W}_L = ((\mathbf{w}')^\top, -(\mathbf{w}')^\top)$ , where each entry of  $\mathbf{w}'$  is generated independently from  $\mathcal{N}(0, 2/m)$ .

**A.1 Verification of Assumption 5.1** For verifying Assumption 5.1, we define that

$$(A.2) \quad g(m) = \max_{t,k} \frac{\|\tilde{\mathbf{f}}_{\text{repr}}(\mathbf{x}_{t,k}; \mathbf{w}_{t,\text{repr}}) - \nabla_{\mathbf{w}} f(\mathbf{x}_{t,k}; \mathbf{w}_t)\|_2}{\sqrt{m}}.$$

Note that by the squeeze theorem if  $\lim_{m \rightarrow +\infty} g(m) = 0$ , then Assumption 1 holds. Therefore, we report  $g(m)$  as well as the corresponding  $m$  in Table 3, where the experiment setting is the same as Section 6.1 and the experiments are repeated for 8 times independently. It shows that  $g(m)$  empirically converges to 0, which plays as a rationale behind Assumption 5.1 from the empirical perspective.

**A.2 Proof of Theorem 5.1** For such a neural network mentioned above, we have that  $\mathbf{f}_{\text{repr}}(\mathbf{x}; \mathbf{w}_{\text{repr}}) = \sqrt{m} \text{ReLU}(\mathbf{f}_{L-1})$  based on Equation 4.4 where  $\mathbf{w}_{\text{repr}} = (\text{vec}(\mathbf{W}_1); \dots; \text{vec}(\mathbf{W}_{L-1})) \in \mathbb{R}^{p-m}$ . For convenience of analysis, we denote  $\lambda = m\gamma$  and  $\nu = \tau^2\gamma$  in Algorithm 2.

We next present the assumptions and definitions used in the proof. Basically, Assumption A.1 requires the reward function is bounded, which is widely adopted in contextual bandit literature [31, 4]. Assumption A.2 is about the contexts and is the same as Assumption 3.4 in [31]. Definition A.1 defines constants  $\sigma$ ,  $s$  that are used a lot in the following proof. Definition A.2 defines the error term  $\epsilon$  and the probability events that are used for our high probability regret bound. Definition A.3 defines the set of saturated arms which are less likely to be selected.

**ASSUMPTION A.1.** There exists an unknown reward function  $h$  such that for any  $1 \leq t \leq T$  and  $1 \leq k \leq K$ ,

$$(A.3) \quad r_{t,k} = h(\mathbf{x}_{t,k}) + \xi_{t,k},$$

where  $|h(\mathbf{x}_{t,k})| \leq 1$ , and  $\{\xi_{t,k}\}$  forms an  $R$ -sub-Gaussian martingale difference sequence with constant  $R > 0$ , i.e.,  $\mathbb{E}[\exp(c\xi_{t,k}) | \xi_{1:t-1,k}, \mathbf{x}_{1:t,k}] \leq \exp(c^2 R^2)$  for all  $c \in \mathbb{R}$ .

**ASSUMPTION A.2.** There exists  $\gamma_0 > 0$ , such that  $\mathbf{H} \succeq \gamma_0 \mathbf{I}$ , where  $\mathbf{H}$  is the neural tangent kernel matrix on the context set, defined the same with [31, 16]. In addition, for any  $t \in [T]$ ,  $k \in [K]$ ,  $\|\mathbf{x}_{t,k}\|_2 = 1$  and  $[\mathbf{x}_{t,k}]_j = [\mathbf{x}_{t,k}]_{j+d/2}$ .

**DEFINITION A.1.** Define  $\sigma_{t,k}$ ,  $s_{t,k}$  and  $\hat{s}_{t,k}$  for  $t \in [T]$ ,  $k \in [K]$  as

$$(A.4) \quad \sigma_{t,k} = \sqrt{\frac{\gamma}{m} \|\mathbf{f}_{\text{repr}}(\mathbf{x}_{t,k}; \mathbf{w}_{t,\text{repr}})\|_{\mathbf{Z}_{t-1}^{-1}}^2},$$

$$(A.5) \quad s_{t,k} = \sqrt{\frac{\gamma}{m} \|\nabla_{\mathbf{w}} f(\mathbf{x}_{t,k}; \mathbf{w}_t)\|_{\mathbf{U}_t^{-1}}^2},$$

$$(A.6) \quad \hat{s}_{t,k} = s_{t,k} - \sigma_{t,k},$$

where  $\mathbf{U}_t = \gamma \mathbf{I} + \sum_{i=1}^{t-1} \nabla_{\mathbf{w}} f(\mathbf{x}_{i,a_i}; \mathbf{w}_{i+1}) \nabla_{\mathbf{w}} f(\mathbf{x}_{i,a_i}; \mathbf{w}_{i+1})^\top / m$ .

**DEFINITION A.2.** If we denote

$$(A.7) \quad \begin{aligned} \epsilon_1(m) &= C_{\epsilon,1}(1 - \eta m \gamma)^{C_{\epsilon,2}} \sqrt{TL/\gamma}, \quad \epsilon_2(m) = \mathcal{O}(\sqrt{\log m m^{-1/6}}), \\ \text{where } C_{\epsilon,1} \text{ and } C_{\epsilon,2} \text{ are positive constants, we could define events } \mathcal{E}_t^\sigma \text{ and } \mathcal{E}_t^\mu \text{ as} \end{aligned}$$

$$(A.8) \quad \mathcal{E}_t^\sigma = \left\{ \omega \in \mathcal{F}_{t+1} : \forall k \in [K], \quad \left| \hat{\theta}_{t,\mathbf{x}_{t,k}} - f(\mathbf{x}_{t,k}; \mathbf{w}_t) \right| \leq c_t \tau \sigma_{t,k} \right\},$$

$$(A.9) \quad \mathcal{E}_t^\mu = \left\{ \omega \in \mathcal{F}_t : \forall k \in [K], \quad |f(\mathbf{x}_{t,k}; \mathbf{w}_t) - h(\mathbf{x}_{t,k})| \leq \tau s_{t,k} + \epsilon(m) \right\}$$

where  $c_t = \sqrt{4 \log t + 2 \log K}$ ,  $\epsilon(m) = \epsilon_1(m) + \epsilon_2(m)$ , and  $\mathcal{F}_t := \{a_i, r_{i,a_i}, \mathbf{x}_{i,k}, k \in [K], i \in [t]\}$  denotes filtration (historical information).

**DEFINITION A.3.** Define the set of saturated arms  $S_t$  as follows

$$(A.10) \quad \begin{aligned} S_t &= \{k \mid k \in [K], h(\mathbf{x}_{t,a_k^*}) - h(\mathbf{x}_{t,k}) \\ &\geq (1 + c_t) \tau \sigma_{t,k} + 2\tau \hat{s}_{t,k} + 2\epsilon(m)\}. \end{aligned}$$

We next introduce the conditions and lemmas in our proof. Basically, Condition A.1 provides the requirement of the width  $m$  for the neural network. Based on Condition A.1, Lemma A.1 and Lemma A.2 bound the probability of the events defined above; Lemma A.3 bounds the regret of each single round.

**CONDITION A.1.** The network width  $m$  satisfies

$$(A.11) \quad m \geq C_m \max\{\sqrt{\gamma} L^{-3/2} [\log(TKL^2/\delta)]^{3/2},$$

$$(A.12) \quad \begin{aligned} T^6 K^6 L^6 \log(TKL/\delta) \max\{\gamma_0^{-4}, 1\}, \\ m[\log m]^{-3} \geq C_m T L^{12} \gamma^{-1} + C_m T^7 \gamma^{-8} L^{18} (\gamma + LT)^6 \\ + C_m L^{21} T^7 \gamma^{-7} (1 + \sqrt{T/\gamma})^6, \end{aligned}$$

where  $C_m$  is a positive constant and  $L$  is the network depth.

**LEMMA A.1.** (Lemma 4.3 in [31]). Suppose the network width  $m$  satisfies Condition A.1, and set  $\eta = C_\eta(m\gamma + mLT)^{-1}$ . Then, it holds that  $\Pr(\forall t \in [T], \mathcal{E}_t^\mu) \geq 1 - \delta$ , where  $C_\eta$  is a positive constant.

**LEMMA A.2.** For all  $t \in [T]$ , the probability for the chosen arm is not in saturated arm set satisfies

$$(A.13) \quad \Pr(a_t \notin S_t \mid \mathcal{F}_t, \mathcal{E}_t^\mu) \geq \frac{1}{4e\sqrt{\pi}} - \frac{1}{t^2}.$$

**LEMMA A.3.** For all  $t \in [T]$ , we have that

$$(A.14) \quad \begin{aligned} &\mathbb{E}[h(\mathbf{x}_{t,a_t^*}) - h(\mathbf{x}_{t,a_t}) \mid \mathcal{F}_t, \mathcal{E}_t^\mu] \\ &\leq \min\{\mathbb{E}[C_1 \tau s_{t,a_t} + C_2 \tau |\sigma_{t,a_t} - s_{t,a_t}| \mid \mathcal{F}_t, \mathcal{E}_t^\mu], 2\} \\ &\quad + 4\epsilon(m) + 2t^{-2}. \end{aligned}$$

where  $C_1$  and  $C_2$  are some positive constants.

Next, we introduce the core lemma used in the proof, Lemma A.4 which decomposes the regret of the algorithm mainly into two parts, i.e.,  $\epsilon_f$  and  $\min\{s_{t,a_t}, 1\}$ . Lemma A.5 and Lemma A.6 bound these parts, respectively.

**LEMMA A.4.** Suppose the network width  $m$  satisfies Condition A.1, and set  $\eta = C_\eta(m\gamma + mLT)^{-1}$ . Then with the probability at least  $1 - \delta$ , it holds that

$$(A.15) \quad \begin{aligned} &\sum_{t=1}^T \left( h(\mathbf{x}_{t,a_t^*}) - h(\mathbf{x}_{t,a_t}) \right) \mathbf{1}(\mathcal{E}_t^\mu) \\ &\leq 4T\epsilon(m) + \pi^2/3 + C_2 \tau T \epsilon_f(m) + C_1 \tau C_s \sqrt{L} \sum_{t=1}^T \min\{s_{t,a_t}, 1\} \\ &\quad + \left( 4 + C_1 \tau C_s \sqrt{L} + C_2 \tau \epsilon_f(m) + 4\epsilon(m) \right) \sqrt{2 \log(1/\delta) T}, \end{aligned}$$

where  $C_\eta, C_1, C_2, C_s$  are some positive constants,  $\epsilon_f(m) = \max_{t \in [T]} |\sigma_{t,a_t} - s_{t,a_t}|$ , and  $\mathbb{1}$  denotes the indicator function for  $\mathcal{E}_t^\mu$ .

LEMMA A.5. (Lemma 4.8 in [34]). Suppose the network width  $m$  satisfies Condition [A.1] and set  $\eta = C_\eta(m\gamma + mLT)^{-1}$ . Then with the probability at least  $1 - \delta$ , it holds that

$$(A.16) \quad \sum_{t=1}^T \min \{s_{t,a_t}, 1\} \leq \sqrt{2\gamma T(C_3 \log(1 + TK) + 1)} + C_4 T^{13/6} \sqrt{\log mm}^{-1/6} \gamma^{-2/3} L^{9/2},$$

where  $C_\eta, C_3, C_4$  are some positive constants.

LEMMA A.6. Under Assumption [5.1], suppose the network width  $m$  satisfies Condition [A.1] and set  $\eta = C_\eta(m\gamma + mLT)^{-1}$ . Then with the probability at least  $1 - \delta$ , it holds that

$$(A.17) \quad \epsilon_f(m) \leq \gamma^{-1/2} m^{-1/2} (1 + C_s^2 LT) \zeta_T$$

where  $C_\eta, C_s$  are some positive constants, and  $\zeta_T = o(\sqrt{m})$ .

The proof of these lemmas are postponed to Appendix [A.3] for the clarification purpose. With all the above lemmas, we are ready to prove Theorem [5.1]

Proof. By Lemma [A.4] with the probability at least  $1 - \delta$ , we have

$$(A.18) \quad \begin{aligned} R_T &= \sum_{t=1}^T \left( h(\mathbf{x}_{t,a_t^*}) - h(\mathbf{x}_{t,a_t}) \right) \mathbb{1}(\mathcal{E}_t^\mu) \\ &\leq 4T\epsilon(m) + \pi^2/3 + C_2\tau T\epsilon_f(m) + C_1\tau C_s\sqrt{L} \sum_{t=1}^T \min \{s_{t,a_t}, 1\} \\ &\quad + \left( 4 + C_1\tau C_s\sqrt{L} + C_2\tau\epsilon_f(m) + 4\epsilon(m) \right) \sqrt{2\log(1/\delta)T} \\ &\leq C_1\tau C_s\sqrt{L} \left( \sqrt{2\gamma T(C_3 \log(1 + TK) + 1)} \right. \\ &\quad \left. + C_2T^{13/6} \sqrt{\log mm}^{-1/6} \gamma^{-2/3} L^{9/2} \right) \\ &\quad + \pi^2/3 + \left( 4 + C_1\tau C_s\sqrt{L} \right) \sqrt{2\log(1/\delta)T} \\ &\quad + \left( 4\epsilon(m) + C_2\tau\epsilon_f(m) \right) \left( T + \sqrt{2\log(1/\delta)T} \right) \\ &\leq C_1\tau C_s\sqrt{L} \sqrt{2\gamma T(C_3 \log(1 + TK) + 1)} \\ &\quad + \left( 4 + C_1\tau C_s\sqrt{L} \right) \sqrt{2\log(1/\delta)T} + \pi^2/3 \\ &\quad + C_1\tau C_s\sqrt{L} C_2T^{13/6} \sqrt{\log mm}^{-1/6} \gamma^{-2/3} L^{9/2} \\ &\quad + 4\epsilon_2(m) \left( T + \sqrt{2\log(1/\delta)T} \right) \\ &\quad + 4C_{\epsilon,1}(1 - \eta m\gamma)^{C_{\epsilon,2}} \sqrt{TL/\gamma} \left( T + \sqrt{2\log(1/\delta)T} \right) \\ &\quad + C_2\tau\gamma^{-1/2} m^{-1/2} (1 + C_s^2 LT) \zeta_T \left( T + \sqrt{2\log(1/\delta)T} \right), \end{aligned}$$

where the first equality comes from the definition of  $R_T$  and  $\mathcal{E}_t^\mu$ , the first inequality comes from Lemma [A.4] the second inequality comes from Lemma [A.5] and the last inequality is from Lemma [A.6]

By setting

$$(A.19) \quad \eta = C_4(m\gamma + mLT)^{-1},$$

$$(A.20) \quad J = (1 + LT/\gamma) \left( \log(24C_{\epsilon,1}) + \log(T^3 L \gamma^{-1} \log(1/\delta)) \right) / C_4,$$

we have

$$(A.21) \quad C_{\epsilon,1}(1 - \eta m\gamma)^{C_{\epsilon,2}} \sqrt{TL/\gamma} (4T + \sqrt{2\log(1/\delta)T}) \leq 1.$$

Based on  $\zeta_T = o(\sqrt{m})$  and  $\epsilon_2(m) = \mathcal{O}(\sqrt{\log mm}^{-1/6})$ , we then can choose  $m$  such that the following three inequalities hold:

$$(A.22) \quad C_1 C_s C_2 T^{13/6} \tau \sqrt{\log mm}^{-1/6} \gamma^{-2/3} L^5 \leq 1,$$

$$(A.23) \quad 4\epsilon_2(m) \left( T + \sqrt{2\log(1/\delta)T} \right) \leq 1,$$

$$(A.24) \quad C_2\tau\gamma^{-1/2} m^{-1/2} (1 + C_s^2 LT) \zeta_T \left( T + \sqrt{2\log(1/\delta)T} \right) \leq 1.$$

We now have the bound:

$$(A.25) \quad \begin{aligned} R_T &\leq C_1\tau C_s\sqrt{L} \sqrt{2\gamma T(C_3 \log(1 + TK) + 1)} \\ &\quad + \left( 4 + C_1\tau C_s\sqrt{L} \right) \sqrt{2\log(1/\delta)T} + 8, \end{aligned}$$

where  $C_1, C_s, C_3$  are positive constants. Taking union bound over Lemmas [A.3], [A.4] and [A.5], the above inequality holds with the probability  $1 - 3\delta$ .

By replacing  $\delta$  with  $3\delta$ , it holds with the probability  $1 - \delta$  that

$$(A.26) \quad R_T = \mathcal{O}(\sqrt{T}),$$

which completes the proof.  $\square$

### A.3 Proof of Lemmas in Section [A.2]

**A.3.1 Proof of Lemma [A.2]** We first give the following two lemmas.

LEMMA A.7. (Lemma 4.2 in [34]). For any  $t \in [T]$ ,  $\Pr(\mathcal{E}_t^\sigma | \mathcal{F}_t) \geq 1 - t^{-2}$ .

LEMMA A.8. (Gaussian anti-concentration). For a Gaussian random variable  $X$  with mean  $\mu$  and standard deviation  $\sigma$ , for any  $\beta > 0$ ,

$$(A.27) \quad \Pr\left(\frac{X - \mu}{\sigma} > \beta\right) \geq \frac{\exp(-\beta^2)}{4\sqrt{\pi}\beta}.$$

Then we start to prove our Lemma [A.2]

Proof. Since  $\hat{\theta}_{t,\mathbf{x}_{t,k}} \sim \mathcal{N}(f(\mathbf{x}_{t,k}; \mathbf{w}_t), \tau^2 \sigma_{t,k}^2)$  conditioned on  $\mathcal{F}_t$ , we have

$$(A.28) \quad \begin{aligned} &\Pr(\hat{\theta}_{t,\mathbf{x}_{t,k}} + \tau \hat{s}_{t,k} + \epsilon(m) > h(\mathbf{x}_{t,k}) | \mathcal{F}_t, \mathcal{E}_t^\mu) \\ &= \Pr\left(\frac{\hat{\theta}_{t,\mathbf{x}_{t,k}} - f(\mathbf{x}_{t,k}; \mathbf{w}_t) + \tau \hat{s}_{t,k} + \epsilon(m)}{\tau \sigma_{t,k}}\right) \\ &> \frac{h(\mathbf{x}_{t,k}) - f(\mathbf{x}_{t,k}; \mathbf{w}_t)}{\tau \sigma_{t,k}} | \mathcal{F}_t, \mathcal{E}_t^\mu) \\ &\geq \Pr\left(\frac{\hat{\theta}_{t,\mathbf{x}_{t,k}} - f(\mathbf{x}_{t,k}; \mathbf{w}_t) + \tau \hat{s}_{t,k}}{\tau \sigma_{t,k}}\right) \\ &> \frac{|h(\mathbf{x}_{t,k}) - f(\mathbf{x}_{t,k}; \mathbf{w}_t)| - \epsilon(m)}{\tau \sigma_{t,k}} | \mathcal{F}_t, \mathcal{E}_t^\mu) \\ &\geq \Pr\left(\frac{\hat{\theta}_{t,\mathbf{x}_{t,k}} - f(\mathbf{x}_{t,k}; \mathbf{w}_t) + \tau \hat{s}_{t,k}}{\tau \sigma_{t,k}} > 1 | \mathcal{F}_t, \mathcal{E}_t^\mu\right) \geq (4e\sqrt{\pi})^{-1}, \end{aligned}$$

where the first inequality is due to  $|x| \geq x$ , and the second inequality is from the definition of  $\mathcal{E}_t^\mu$ .

Consider the following two events at round  $t$ :

$$(A.29) \quad \mathcal{A} = \left\{ \forall k \in S_t, \hat{\theta}_{t,\mathbf{x}_{t,k}} < \hat{\theta}_{t,\mathbf{x}_{t,a_t^*}} | \mathcal{F}_t, \mathcal{E}_t^\mu \right\}$$

$$(A.30) \quad \mathcal{B} = \{a_t \notin S_t | \mathcal{F}_t, \mathcal{E}_t^\mu\}$$

Clearly,  $\mathcal{A}$  implies  $\mathcal{B}$ , since  $a_t = \arg\max_k \hat{\theta}_{t,\mathbf{x}_{t,k}}$ . Therefore,

$$(A.31) \quad \Pr(a_t \notin S_t | \mathcal{F}_t, \mathcal{E}_t^\mu) \geq \Pr(\forall k \in S_t, \hat{\theta}_{t,\mathbf{x}_{t,k}} < \hat{\theta}_{t,\mathbf{x}_{t,a_t^*}} | \mathcal{F}_t, \mathcal{E}_t^\mu).$$

Suppose  $\mathcal{E}_t^\mu$  and  $\mathcal{E}_t^\sigma$  also hold, then it is easy to show that  $\forall k \in [K]$ ,

$$(A.32) \quad \begin{aligned} |h(\mathbf{x}_{t,k}) - \hat{\theta}_{t,\mathbf{x}_{t,k}}| &\leq |h(\mathbf{x}_{t,k}) - f(\mathbf{x}_{t,k}; \mathbf{w}_{t-1})| \\ &\quad + |f(\mathbf{x}_{t,k}; \mathbf{w}_{t-1}) - \hat{\theta}_{t,\mathbf{x}_{t,k}}| \\ &\leq \tau s_{t,k} + \epsilon(m) + c_t \tau \sigma_{t,k} \\ &= (1 + c_t) \tau \sigma_{t,k} + \tau \hat{s}_{t,k} + \epsilon(m), \end{aligned}$$

where the first inequality is due to  $|x_1 + x_2| \leq |x_1| + |x_2|$ , the second inequality is from Definition [A.2], and the last equality is from Definition [A.1]

Hence, for all  $k \in S_t$ , we have that

$$(A.33) \quad h(\mathbf{x}_{t,a_t^*}) - \hat{\theta}_{t,\mathbf{x}_{t,k}} \geq h(\mathbf{x}_{t,a_t^*}) - h(\mathbf{x}_{t,k}) - |h(\mathbf{x}_{t,k}) - \hat{\theta}_{t,\mathbf{x}_{t,k}}| \\ \geq \tau \hat{s}_{t,k} + \epsilon(m),$$

where the first inequality is due to  $|x| \geq x$ , and the second inequality is from Definition [A.3](#) and Equation [A.32](#).

Consider the following event

$$(A.34) \quad \mathcal{C} = \left\{ h(\mathbf{x}_{t,a_t^*}) - \tau \hat{s}_{t,k} - \epsilon(m) < \hat{\theta}_{t,\mathbf{x}_{t,a_t^*}} \mid \mathcal{F}_t, \mathcal{E}_t^\mu \right\}.$$

Since  $\mathcal{E}_t^\mu$  holds,  $\mathcal{E}_t^\sigma$  then implies [A.33](#) as shown above. We then have that if  $\mathcal{C}$  and  $\mathcal{E}_t^\sigma$  hold, then  $\mathcal{A}$  holds, i.e.  $\mathcal{E}_t^\sigma \cap \mathcal{C} \subseteq \mathcal{A}$ . Taking union with  $\bar{\mathcal{E}}_t^\sigma$ , we have that  $\mathcal{C} = \bar{\mathcal{E}}_t^\sigma \cup \mathcal{E}_t^\sigma \cap \mathcal{C} \subseteq \mathcal{A} \cup \bar{\mathcal{E}}_t^\sigma$ , which implies

$$(A.35) \quad \Pr(\mathcal{A}) + \Pr(\bar{\mathcal{E}}_t^\sigma) \geq \Pr(\mathcal{C}).$$

Then, we have that

$$(A.36) \quad \Pr(a_t \notin S_t \mid \mathcal{F}_t, \mathcal{E}_t^\mu) \\ \geq \Pr(\forall k \in S_t, \hat{\theta}_{t,\mathbf{x}_{t,k}} < \hat{\theta}_{t,\mathbf{x}_{t,a_t^*}} \mid \mathcal{F}_t, \mathcal{E}_t^\mu) \\ \geq \Pr(\hat{\theta}_{t,\mathbf{x}_{t,a_t^*}} + \tau \hat{s}_{t,k} + \epsilon(m) > h(\mathbf{x}_{t,a_t^*}) \mid \mathcal{F}_t, \mathcal{E}_t^\mu) \\ - \Pr(\bar{\mathcal{E}}_t^\sigma \mid \mathcal{F}_t, \mathcal{E}_t^\mu) \\ \geq \frac{1}{4e\sqrt{\pi}} - \frac{1}{t^2},$$

where the first equality is due to [A.31](#), the first inequality is due to [A.35](#) and [A.34](#), and the second inequality is from Lemmas [A.7](#) and [A.28](#).  $\square$

### A.3.2 Proof of Lemma [A.3](#)

*Proof.* Recall that given  $\mathcal{F}_t$  and  $\mathcal{E}_t^\mu$ , the only randomness comes from sampling  $\hat{\theta}_{t,\mathbf{x}_{t,k}}$  for  $k \in [K]$ .

We denote  $\bar{k}_t$  as

$$(A.37) \quad \bar{k}_t = \underset{k \notin S_t}{\operatorname{argmin}} \quad 2(1+c_t)\tau\sigma_{t,k} + 3\tau\hat{s}_{t,k}.$$

Then we have

$$(A.38) \quad \mathbb{E}[2(1+c_t)\tau\sigma_{t,a_t} + 3\tau\hat{s}_{t,a_t} \mid \mathcal{F}_t, \mathcal{E}_t^\mu] \\ \geq \mathbb{E}[2(1+c_t)\tau\sigma_{t,a_t} + 3\tau\hat{s}_{t,a_t} \mid \mathcal{F}_t, \mathcal{E}_t^\mu, a_t \notin S_t] \Pr(a_t \notin S_t \mid \mathcal{F}_t, \mathcal{E}_t^\mu) \\ \geq (2(1+c_t)\tau\sigma_{t,\bar{k}_t} + 3\tau\hat{s}_{t,\bar{k}_t}) \left( \frac{1}{4e\sqrt{\pi}} - \frac{1}{t^2} \right),$$

where the first inequality is due to the property of conditional expectation, and the second inequality is from Lemma [A.2](#) and [A.37](#).

If both  $\mathcal{E}_t^\sigma$  and  $\mathcal{E}_t^\mu$  hold, then

$$(A.39) \quad \forall k \in [K], |h(\mathbf{x}_{t,k}) - \hat{\theta}_{t,\mathbf{x}_{t,k}}| \leq \epsilon(m) + (1+c_t)\tau\sigma_{t,k} + \tau\hat{s}_{t,k}$$

as proved in [A.32](#). Thus,

$$(A.40) \quad h(\mathbf{x}_{t,a_t^*}) - h(\mathbf{x}_{t,a_t}) \\ = h(\mathbf{x}_{t,a_t^*}) - h(\mathbf{x}_{t,\bar{k}_t}) + h(\mathbf{x}_{t,\bar{k}_t}) - h(\mathbf{x}_{t,a_t}) \\ \leq (1+c_t)\tau\sigma_{t,\bar{k}_t} + 2\tau\hat{s}_{t,\bar{k}_t} + 2\epsilon(m) + h(\mathbf{x}_{t,\bar{k}_t}) \\ - \hat{\theta}_{t,\mathbf{x}_{t,\bar{k}_t}} - h(\mathbf{x}_{t,a_t}) + \hat{\theta}_{t,\mathbf{x}_{t,a_t}} + \hat{\theta}_{t,\mathbf{x}_{t,\bar{k}_t}} - \hat{\theta}_{t,\mathbf{x}_{t,a_t}} \\ \leq (1+c_t)\tau\sigma_{t,\bar{k}_t} + 2\tau\hat{s}_{t,\bar{k}_t} + 2\epsilon(m) \\ + |h(\mathbf{x}_{t,\bar{k}_t}) - \hat{\theta}_{t,\mathbf{x}_{t,\bar{k}_t}}| + |h(\mathbf{x}_{t,a_t}) - \hat{\theta}_{t,\mathbf{x}_{t,a_t}}| \\ \leq (1+c_t)\tau\sigma_{t,\bar{k}_t} + 2\tau\hat{s}_{t,\bar{k}_t} + 2\epsilon(m) + \epsilon(m) + (1+c_t)\tau\sigma_{t,\bar{k}_t} + \tau\hat{s}_{t,\bar{k}_t} \\ + \epsilon(m) + (1+c_t)\tau\sigma_{t,a_t} + \tau\hat{s}_{t,a_t} \\ = (1+c_t)\tau(2\sigma_{t,\bar{k}_t} + \sigma_{t,a_t}) + 3\tau\hat{s}_{t,\bar{k}_t} + \tau\hat{s}_{t,a_t} + 4\epsilon(m),$$

where the first inequality is due to Definition [A.3](#) and  $\bar{k}_t \notin S_t$ , the second inequality is due to  $|x| \geq x$  and  $\hat{\theta}_{t,\mathbf{x}_{t,\bar{k}_t}} \leq \hat{\theta}_{t,\mathbf{x}_{t,a_t}}$ , and the last inequality is from [A.39](#).

Since a trivial bound on  $h(\mathbf{x}_{t,a_t^*}) - h(\mathbf{x}_{t,a_t})$  could be derived

by

$$(A.41) \quad h(\mathbf{x}_{t,a_t^*}) - h(\mathbf{x}_{t,a_t}) \leq |h(\mathbf{x}_{t,a_t^*})| + |h(\mathbf{x}_{t,a_t})| \leq 2,$$

then we have

$$(A.42) \quad \mathbb{E}[h(\mathbf{x}_{t,a_t^*}) - h(\mathbf{x}_{t,a_t}) \mid \mathcal{F}_t, \mathcal{E}_t^\mu] \\ = \mathbb{E}[h(\mathbf{x}_{t,a_t^*}) - h(\mathbf{x}_{t,a_t}) \mid \mathcal{F}_t, \mathcal{E}_t^\mu, \mathcal{E}_t^\sigma] \Pr(\mathcal{E}_t^\sigma) \\ + \mathbb{E}[h(\mathbf{x}_{t,a_t^*}) - h(\mathbf{x}_{t,a_t}) \mid \mathcal{F}_t, \mathcal{E}_t^\mu, \bar{\mathcal{E}}_t^\sigma] \Pr(\bar{\mathcal{E}}_t^\sigma) \\ \leq (1+c_t)\tau(2\sigma_{t,\bar{k}_t} + \mathbb{E}[\sigma_{t,a_t} \mid \mathcal{F}_t, \mathcal{E}_t^\mu]) \\ + 4\epsilon(m) + \frac{2}{t^2} + 3\tau\hat{s}_{t,\bar{k}_t} + \tau\mathbb{E}[\hat{s}_{t,a_t} \mid \mathcal{F}_t, \mathcal{E}_t^\mu] \\ = 2(1+c_t)\tau\sigma_{t,\bar{k}_t} + 3\tau\hat{s}_{t,\bar{k}_t} + 4\epsilon(m) \\ + \frac{2}{t^2} + \mathbb{E}[(1+c_t)\tau\sigma_{t,a_t} + \tau\hat{s}_{t,a_t} \mid \mathcal{F}_t, \mathcal{E}_t^\mu] \\ \leq \frac{\mathbb{E}[2(1+c_t)\tau\sigma_{t,a_t} + 3\tau\hat{s}_{t,a_t} \mid \mathcal{F}_t, \mathcal{E}_t^\mu]}{\frac{1}{4e\sqrt{\pi}} - \frac{1}{t^2}} \\ + 4\epsilon(m) + \frac{2}{t^2} + \mathbb{E}[(1+c_t)\tau\sigma_{t,a_t} + \tau\hat{s}_{t,a_t} \mid \mathcal{F}_t, \mathcal{E}_t^\mu]$$

$\leq \mathbb{E}[44e\sqrt{\pi}(1+c_t)\tau\sigma_{t,a_t} + 64e\sqrt{\pi}\tau\hat{s}_{t,a_t} \mid \mathcal{F}_t, \mathcal{E}_t^\mu] + 4\epsilon(m) + 2t^{-2}$  where the first equality is due to the property of conditional expectation, the first inequality uses the bound provide in [A.40](#) and the trivial bound of  $h(\mathbf{x}_{t,a_t^*}) - h(\mathbf{x}_{t,a_t})$  for the second term plus Lemma [A.7](#), the second inequality uses the bound provided in [A.38](#), the last inequality is directly calculated by  $1 \leq 4e\sqrt{\pi}$  and

$$(A.43) \quad \frac{1}{\frac{1}{4e\sqrt{\pi}} - \frac{1}{t^2}} \leq 20e\sqrt{\pi},$$

which holds since LHS is negative when  $t \leq 4$  and when  $t = 5$ , the LHS reach its maximum as  $\approx 84.11 < 96.36 \approx \text{RHS}$ .

Noticing that  $|h(\mathbf{x})| \leq 1$ , it is trivial to further extend the bound as

$$(A.44) \quad \mathbb{E}[h(\mathbf{x}_{t,a_t^*}) - h(\mathbf{x}_{t,a_t}) \mid \mathcal{F}_t, \mathcal{E}_t^\mu] \\ \leq \min\{\mathbb{E}[44e\sqrt{\pi}(1+c_t)\tau\sigma_{t,a_t} + 64e\sqrt{\pi}\tau\hat{s}_{t,a_t} \mid \mathcal{F}_t, \mathcal{E}_t^\mu], 2\} \\ + 4\epsilon(m) + 2t^{-2} \\ = \min\{\mathbb{E}[44e\sqrt{\pi}(1+c_t)\tau\sigma_{t,a_t} \\ + (44e\sqrt{\pi}(1+c_t) - 64e\sqrt{\pi})\tau(\sigma_{t,a_t} - s_{t,a_t}) \\ \mid \mathcal{F}_t, \mathcal{E}_t^\mu], 2\} + 4\epsilon(m) + 2t^{-2} \\ \leq \min\{\mathbb{E}[44e\sqrt{\pi}(1+c_T)\tau\sigma_{t,a_t} \\ + (44e\sqrt{\pi}(1+c_T) - 64e\sqrt{\pi})\tau|\sigma_{t,a_t} - s_{t,a_t}| \\ \mid \mathcal{F}_t, \mathcal{E}_t^\mu], 2\} + 4\epsilon(m) + 2t^{-2}.$$

Let denote  $C_1 = 44e\sqrt{\pi}(1+c_T)$  and  $C_2 = 44e\sqrt{\pi}(1+c_T) - 64e\sqrt{\pi}$ . We have that

$$(A.45) \quad \mathbb{E}[h(\mathbf{x}_{t,a_t^*}) - h(\mathbf{x}_{t,a_t}) \mid \mathcal{F}_t, \mathcal{E}_t^\mu] \\ \leq \min\{\mathbb{E}[C_1\tau\sigma_{t,a_t} + C_2\tau|\sigma_{t,a_t} - s_{t,a_t}| \mid \mathcal{F}_t, \mathcal{E}_t^\mu], 2\} + 4\epsilon(m) + 2t^{-2}. \\ \square$$

### A.3.3 Proof of Lemma [A.4](#)

We first introduce two lemmas to support our proof.

LEMMA A.9. (Lemma B.9 in [\[37\]](#)). For any  $t \in [T]$ ,  $k \in [K]$  and  $\delta \in (0, 1)$ , if the network width  $m$  satisfies Condition [A.1](#), we have, with probability at least  $1 - \delta$ , that

$$(A.46) \quad s_{t,k} \leq C_s \sqrt{L},$$

where  $C_s$  is a positive constant.

LEMMA A.10. (Azuma-Hoeffding Inequality for Super Martingale). If a super-martingale  $Y_t$ , corresponding to filtration  $\mathcal{F}_t$  satisfies that  $|Y_t - Y_{t-1}| \leq B_t$ , then for any  $\delta \in (0, 1)$ , we have

$$(A.47) \quad Y_t - Y_0 \leq \sqrt{2 \log(1/\delta) \sum_{i=1}^t B_i^2}.$$

Then the proof for Lemma A.4 is as follows:

*Proof.* By setting  $\tau$  such that  $C_1\tau \geq 2$  and using Lemma A.3 we have

$$(A.48) \quad \begin{aligned} & \mathbb{E} \left[ h(\mathbf{x}_{t,a_t^*}) - h(\mathbf{x}_{t,a_t}) \mid \mathcal{F}_t, \mathcal{E}_t^\mu \right] \\ & \leq \min \{ \mathbb{E} [C_1\tau s_{t,a_t} \mid \mathcal{F}_t, \mathcal{E}_t^\mu], 2 \} \\ & \quad + \min \{ \mathbb{E} [C_2\tau |\sigma_{t,a_t} - s_{t,a_t}| \mid \mathcal{F}_t, \mathcal{E}_t^\mu], 2 \} \\ & \quad + 4\epsilon(m) + 2t^{-2} \\ & \leq C_1\tau \min \{ \mathbb{E} [s_{t,a_t} \mid \mathcal{F}_t, \mathcal{E}_t^\mu], 1 \} + \mathbb{E} [C_2\tau |\sigma_{t,a_t} - s_{t,a_t}| \mid \mathcal{F}_t, \mathcal{E}_t^\mu] \\ & \quad + 4\epsilon(m) + 2t^{-2}, \end{aligned}$$

where the first inequality is due to  $\min\{|x|+|y|, 2\} \leq \min\{|x|, 2\} + \min\{|y|, 2\}$ , and the last inequality is due to  $C_1\tau \geq 2$  and  $\min\{x, 2\} \leq x$ .

Based on Lemma A.9, with probability at least  $1 - \delta$ , we have

$$(A.49) \quad \begin{aligned} & \mathbb{E} \left[ h(\mathbf{x}_{t,a_t^*}) - h(\mathbf{x}_{t,a_t}) \mid \mathcal{F}_t, \mathcal{E}_t^\mu \right] \\ & \leq C_1\tau C_s \sqrt{L} \mathbb{E} [\min \{s_{t,a_t}, 1\} \mid \mathcal{F}_t, \mathcal{E}_t^\mu] \\ & \quad + \mathbb{E} [C_2\tau |\sigma_{t,a_t} - s_{t,a_t}| \mid \mathcal{F}_t, \mathcal{E}_t^\mu] \\ & \quad + 4\epsilon(m) + 2t^{-2}, \end{aligned}$$

where  $C_s$  is a positive constant.

We define  $\Delta_t := \left( h(\mathbf{x}_{t,a_t^*}) - h(\mathbf{x}_{t,a_t}) \right) \mathbf{1}(\mathcal{E}_t^\mu)$ , and

$$(A.50) \quad \begin{aligned} X_t &:= \Delta_t - \left( C_1\tau C_s \sqrt{L} \min \{s_{t,a_t}, 1\} \right. \\ & \quad \left. + C_2\tau |\sigma_{t,a_t} - s_{t,a_t}| + 4\epsilon(m) + 2t^{-2} \right), \\ Y_t &= \sum_{i=1}^t X_i. \end{aligned}$$

Then, with probability at least  $1 - \delta$ , we have

$$(A.51) \quad \begin{aligned} |X_t| &\leq |\Delta_t| + C_1\tau C_s \sqrt{L} \min \{s_{t,a_t}, 1\} \\ &\quad + C_2\tau |\sigma_{t,a_t} - s_{t,a_t}| + 4\epsilon(m) + 2t^{-2} \\ &\leq 2 + 2t^{-2} + C_1\tau C_s \sqrt{L} + C_2\tau |\sigma_{t,a_t} - s_{t,a_t}| + 4\epsilon(m) \\ &\leq 4 + C_1\tau C_s \sqrt{L} + C_2\tau \epsilon_f(m) + 4\epsilon(m) \end{aligned}$$

where the first inequality is due to  $|x+y| \leq |x|+|y|$ , the second inequality uses the fact that  $h \leq 1$ , and the second inequality is due to  $\epsilon_f(m) := \max_{t \in [T]} |\sigma_{t,a_t} - s_{t,a_t}|$ .

From Lemma A.10, with probability at least  $1 - \delta$ , we have

$$(A.52) \quad \begin{aligned} \sum_{t=1}^T X_t &= Y_T - Y_0 \\ &\leq \left( 4 + C_1\tau C_s \sqrt{L} + C_2\tau \epsilon_f(m) + 4\epsilon(m) \right) \sqrt{2 \log(1/\delta)T}, \end{aligned}$$

We then can derive that

$$(A.53) \quad \begin{aligned} \sum_{t=1}^T \Delta_t &\leq \left( 4 + C_1\tau C_s \sqrt{L} + C_2\tau \epsilon_f(m) + 4\epsilon(m) \right) \sqrt{2 \log(1/\delta)T} \\ &\quad + \sum_{t=1}^T \left( C_1\tau C_s \sqrt{L} \min \{s_{t,a_t}, 1\} + C_2\tau |\sigma_{t,a_t} - s_{t,a_t}| \right. \\ &\quad \left. + 4\epsilon(m) + 2t^{-2} \right) \\ &\leq 4T\epsilon(m) + \pi^2/3 + C_2\tau T \epsilon_f(m) \\ &\quad + C_1\tau C_s \sqrt{L} \sum_{t=1}^T \min \{s_{t,a_t}, 1\} \\ &\quad + \left( 4 + C_1\tau C_s \sqrt{L} + C_2\tau \epsilon_f(m) + 4\epsilon(m) \right) \sqrt{2 \log(1/\delta)T}, \end{aligned}$$

where the first inequality is due to (A.50) and (A.52), and the last inequality is due to the definition of  $\epsilon_f(m)$  and  $\sum_{i=1}^\infty t^{-2} = \pi^2/6$ .

□

**A.3.4 Proof of Lemma A.6** We first present our Lemma A.11.

LEMMA A.11. (Lemma B.4 in [37]). Define  $\psi(\mathbf{a}, \mathbf{a}_1, \dots, \mathbf{a}_{t-1})$  as

$$(A.54) \quad \psi(\mathbf{a}, \mathbf{a}_1, \dots, \mathbf{a}_{t-1}) := \sqrt{\mathbf{a}^\top \left( \gamma \mathbf{I} + \sum_{i=1}^{t-1} \mathbf{a}_i \mathbf{a}_i^\top \right)^{-1} \mathbf{a}}.$$

If  $\|\mathbf{a}\|_2 \leq C_{Lip} \sqrt{L}$ , then  $\psi(\mathbf{a}, \mathbf{a}_1, \dots, \mathbf{a}_{t-1})$  is Lipschitz continuous, where  $\|\nabla_{\mathbf{a}} \psi\|_2 \leq 1/\sqrt{\gamma}$  and  $\|\nabla_{\mathbf{a}_i} \psi\|_2 \leq C_{Lip}^2 L/\sqrt{\gamma}$ .

Then with the lemma above, we derive our Lemma A.6.

*Proof.* It is easy to verify that

$$(A.55) \quad \begin{aligned} \sigma_{t,k} &= \psi \left( \frac{\tilde{\mathbf{f}}_{\text{repr}}(\mathbf{x}_{t,k}; \mathbf{w}_{t,\text{repr}})}{\sqrt{m}}, \right. \\ & \quad \left. \frac{\tilde{\mathbf{f}}_{\text{repr}}(\mathbf{x}_{1,a_1}; \mathbf{w}_{2,\text{repr}})}{\sqrt{m}}, \dots, \frac{\tilde{\mathbf{f}}_{\text{repr}}(\mathbf{x}_{t-1,a_{t-1}}; \mathbf{w}_{t,\text{repr}})}{\sqrt{m}} \right) \\ s_{t,k} &= \psi \left( \frac{\nabla_{\mathbf{w}} f(\mathbf{x}_{t,k}; \mathbf{w}_t)}{\sqrt{m}}, \right. \\ & \quad \left. \frac{\nabla_{\mathbf{w}} f(\mathbf{x}_{1,a_1}; \mathbf{w}_2)}{\sqrt{m}}, \dots, \frac{\nabla_{\mathbf{w}} f(\mathbf{x}_{t-1,a_{t-1}}; \mathbf{w}_t)}{\sqrt{m}} \right). \end{aligned}$$

where  $\tilde{\mathbf{f}}_{\text{repr}}$  is expanded from  $\mathbf{f}_{\text{repr}}$  by aligning its elements with  $\text{vec}(\mathbf{W}_L)$  in  $\mathbf{w}_t$  and adding 0 in other positions.

By Lemma A.11, Lemma A.9 and Equation A.55 we have that:

$$(A.56) \quad \begin{aligned} |\sigma_{t,a_t} - s_{t,a_t}| &\leq \frac{1}{\sqrt{\gamma}} \left\| \frac{\tilde{\mathbf{f}}_{\text{repr}}(\mathbf{x}_{t,k}; \mathbf{w}_{t,\text{repr}})}{\sqrt{m}} - \frac{\nabla_{\mathbf{w}} f(\mathbf{x}_{t,k}; \mathbf{w}_t)}{\sqrt{m}} \right\|_2 \\ &\quad + \frac{C_s^2 L}{\sqrt{\gamma}} \sum_{i=1}^{t-1} \left\| \frac{\tilde{\mathbf{f}}_{\text{repr}}(\mathbf{x}_{i,a_i}; \mathbf{w}_{i+1,\text{repr}})}{\sqrt{m}} - \frac{\nabla_{\mathbf{w}} f(\mathbf{x}_{i,a_i}; \mathbf{w}_{i+1})}{\sqrt{m}} \right\|_2. \end{aligned}$$

Then, we have

$$(A.57) \quad \epsilon_f(m) \leq \gamma^{-1/2} m^{-1/2} \left( T^{-1} + C_s^2 L \right) T \zeta_T$$

where  $\zeta_T = \max\{\zeta_{T,1}, \zeta_{T,2}\}$  with

$$(A.58) \quad \zeta_{T,1} = \max_{t \in [T], k \in [K]} \left\| \tilde{\mathbf{f}}_{\text{repr}}(\mathbf{x}_{t,k}; \mathbf{w}_{t,\text{repr}}) - \nabla_{\mathbf{w}} f(\mathbf{x}_{t,k}; \mathbf{w}_t) \right\|_2,$$

$$(A.59) \quad \zeta_{T,2} = \max_{t \in [T]} \left\| \tilde{\mathbf{f}}_{\text{repr}}(\mathbf{x}_{t,a_t}; \mathbf{w}_{t+1,\text{repr}}) - \nabla_{\mathbf{w}} f(\mathbf{x}_{t,a_t}; \mathbf{w}_{t+1}) \right\|_2.$$

By Assumption 5.1 we know that  $\zeta_T = o(\sqrt{m})$ .

□

## B Details of Algorithm Variants

In this section, we introduce two variants of PlugTS in detail, Practical PlugTS (PlugTS-P) and Generalized PlugTS (PlugTS-G).

In PlugTS-P, we do not maintain the historical context information in the matrix  $\mathbf{M}$  explicitly, which further reduces computational overhead. Similar with Algorithm 2, we formally present PlugTS-P in Algorithm 3.

In PlugTS-G, we utilize gradients of parameters for multiple layers instead of only the gradient with respect to parameters of the last layer. Note that  $\nabla_{\text{vec}(\mathbf{W}_L)} f(\mathbf{x}; \mathbf{w}) = \mathbf{f}_{\text{repr}}(\mathbf{x}; \mathbf{w}_{\text{repr}})$  for the neural network defined in (A.1). We then replace  $\nabla_{\text{vec}(\mathbf{W}_L)} f(\mathbf{x}; \mathbf{w})$  with  $\nabla_{\mathbf{w}_{\text{subset}}} f(\mathbf{x}; \mathbf{w})$ , where  $\mathbf{w}_{\text{subset}} = (\text{vec}(\mathbf{W}_{s_1}); \dots; \text{vec}(\mathbf{W}_{s_l}))$ ,  $s_1, \dots, s_l \in [L]$ , i.e.,  $\mathbf{w}_{\text{subset}}$  is the collection vector of parameters for some layers.

PlugTS-G will reduce to PlugTS when  $\mathbf{w}_{\text{subset}} = \text{vec}(\mathbf{W}_L)$ , and reduce to NeuralTS [31] when  $\mathbf{w}_{\text{subset}} = \mathbf{w}$ . We present PlugTS-G in Algorithm 4 with denoting  $\mathbf{f}_{\text{subset}}(\mathbf{x}; \mathbf{w}_{\text{subset}}) := \nabla_{\mathbf{w}_{\text{subset}}} f(\mathbf{x}; \mathbf{w})$ .

We then generalize Assumption 5.1 to fit PlugTS-G.

ASSUMPTION B.1. Define that  $\tilde{\mathbf{f}}_{\text{subset}}$  is expanded from  $\mathbf{f}_{\text{subset}}$  by aligning its elements with  $\mathbf{w}_{t,\text{subset}}$  in  $\mathbf{w}_t$  and padding with 0

**Algorithm 3: PlugTS-P (Practical PlugTS)**

**Input:** Initialization parameter of neural network  $\mathbf{w}_1$ , hyper-parameter for exploration strength  $\nu$ , network width  $m$ , regularization parameter  $\lambda > 0$ , positive definite matrix  $\mathbf{M}_1 = \lambda \mathbf{I}$

```

1 for  $t = 1, \dots, T$  do
2   Receive context vectors  $\{\mathbf{x}_{t,k}\}_{k=1}^K$ ;
3   for  $k = 1, \dots, K$  do
4     Sample  $\hat{\theta}_{t,\mathbf{x}_{t,k}} \sim \mathcal{N}(f(\mathbf{x}_{t,k}; \mathbf{w}_t), \nu \|\mathbf{f}_{\text{repr}}(\mathbf{x}_{t,k}; \mathbf{w}_t, \text{repr})\|_{\mathbf{M}_t^{-1}}^2)$ ;
5   end
6    $a_t \leftarrow \arg \max_{a \in [K]} \hat{\theta}_{t,\mathbf{x}_{t,a}}$ ;
7   Apply  $a_t$  and observe  $r_{t,a_t}$ ;
8    $\mathbf{w}_{t+1} \leftarrow \arg \min_{\mathbf{w}} \sum_{i=1}^t (f(\mathbf{x}_{t,a_t}) - r_{t,a_t})^2 / 2 + \lambda \|\mathbf{w} - \mathbf{w}_1\|^2 / 2$ ;
9    $\mathbf{M}_{t+1} \leftarrow \mathbf{M}_t + \mathbf{I}$ ;
10 end

```

**Algorithm 4: PlugTS-G (Generalized PlugTS)**

**Input:** Initialization parameter of neural network  $\mathbf{w}_1$ , hyper-parameter for exploration strength  $\tau$ , network width  $m$ , regularization parameter  $\gamma > 0$ , positive definite matrix  $\mathbf{Z}_1 = \gamma \mathbf{I}$

```

1 for  $t = 1, \dots, T$  do
2   Receive context vectors  $\{\mathbf{x}_{t,k}\}_{k=1}^K$ ;
3   for  $k = 1, \dots, K$  do
4     Sample  $\hat{\theta}_{t,\mathbf{x}_{t,k}} \sim \mathcal{N}(f(\mathbf{x}_{t,k}; \mathbf{w}_t), \tau^2 \gamma \|\mathbf{f}_{\text{subset}}(\mathbf{x}_{t,k}; \mathbf{w}_t, \text{subset})\|_{\mathbf{Z}_t^{-1}}^2 / m)$ ;
5   end
6    $a_t \leftarrow \arg \max_{a \in [K]} \hat{\theta}_{t,\mathbf{x}_{t,a}}$ ;
7   Apply  $a_t$  and observe  $r_{t,a_t}$ ;
8    $\mathbf{w}_{t+1} \leftarrow \arg \min_{\mathbf{w}} \sum_{i=1}^t (f(\mathbf{x}_{t,a_t}) - r_{t,a_t})^2 / 2 + m \gamma \|\mathbf{w} - \mathbf{w}_1\|^2 / 2$ ;
9    $\mathbf{Z}_{t+1} \leftarrow \mathbf{Z}_t + \mathbf{f}_{\text{subset}}(\mathbf{x}_{t,a_t}; \mathbf{w}_{t+1}, \text{subset}) \mathbf{f}_{\text{subset}}(\mathbf{x}_{t,a_t}; \mathbf{w}_{t+1}, \text{subset})^\top / m$ ;
10 end

```

Table 4: Total reward averaged over 8 repeated experiments.

	PlugTS	PlugTS-G-2	PlugTS-G-3	NeuralTS
Yahoo!R6B	52381.8 $\pm$ 540.2	52346.2 $\pm$ 660.4	52237.0 $\pm$ 602.9	52235.1 $\pm$ 1809.3

for the other elements. We assume it satisfies

$$(B.60) \quad \|\mathbf{f}_{\text{subset}}(\mathbf{x}_{t,k}; \mathbf{w}_t, \text{subset}) - \nabla_{\mathbf{w}} f(\mathbf{x}_{t,k}; \mathbf{w}_t)\| = o(\sqrt{m}).$$

Under Assumption B.1 it is easy to verify that the proof in Appendix A.2 still holds with replacing  $\mathbf{f}_{\text{repr}}(\mathbf{x}; \mathbf{w}_{\text{repr}})$  with  $\mathbf{f}_{\text{subset}}(\mathbf{x}; \mathbf{w}_{\text{subset}})$ . Therefore, PlugTS-G also achieves the  $\mathcal{O}(\sqrt{T})$  regret bound.

Moreover, we have realized two versions of PlugTS-G, named PlugTS-G-2 and PlugTS-G-3, where PlugTS-G-2 uses the gradients of the last two layers and PlugTS-G-3 uses the gradients of the last three layers. The experiment setting are the same as Section 6.2 and the results over 8 repeated experiments on the Yahoo!R6B dataset are reported in Table 4. The results show that there is no significant difference in the performance of the four algorithms, which furthermore verifies our theoretical finding that the full gradient is not necessary for the sublinear regret bound.

**Algorithm 5: Online Recommender System Simulator**

**Input:** randomized impression dataset  $\mathcal{D}$ , initialization dataset  $\mathcal{D}_{\text{init}}$ , selecting strategy  $s : \mathbb{X} \times \mathbb{A} \rightarrow \mathbb{R}$  and corresponding neural network  $f$ , size of partition  $N > 0$

```

1  $\mathcal{D}_{\text{train}} \leftarrow \mathcal{D}_{\text{init}}$ ;
2 Train  $f$  with  $\mathcal{D}_{\text{train}}$ ;
3  $r_{\text{all}} \leftarrow 0$  //  $r_{\text{all}}$ : cumulative reward (number of clicks);
4 for each  $N$ -size log entries partition  $\mathcal{D}_{\text{part}} \subseteq \mathcal{D}$  do
5   Initialize  $\mathcal{D}_{\text{simulator}} \leftarrow \emptyset$ ;
6   //  $\mathbf{x}$ : user feature;  $\mathbf{a}$ : displayed item;  $\mathbf{r}$ : binary label;
7   A: candidate set for  $(\mathbf{x}, \mathbf{a}, \mathbf{r}, \mathbf{A}) \in \mathcal{D}_{\text{part}}$  do
8     if  $\arg \max_{\mathbf{a}' \in \mathbf{A}} s(\mathbf{x}, \mathbf{a}') = \mathbf{a}$  then
9        $\mathcal{D}_{\text{simulator}} \leftarrow \mathcal{D}_{\text{simulator}} \cup \{(\mathbf{x}, \mathbf{a}, \mathbf{r})\}$ ;
10       $r_{\text{all}} \leftarrow r_{\text{all}} + r$ ;
11    end
12  Remove the oldest  $|\mathcal{D}_{\text{simulator}}|$  samples from  $\mathcal{D}_{\text{train}}$ ;
13   $\mathcal{D}_{\text{train}} \leftarrow \mathcal{D}_{\text{train}} \cup \mathcal{D}_{\text{simulator}}$ ;
14  Train  $f$  with  $\mathcal{D}_{\text{train}}$ ;
15 end
Output:  $r_{\text{all}}$ 

```

Table 5: Total regret on classification datasets averaged over 8 repeated experiments.

	Greedy	MCDropout	NeuralTS	PlugTS	PlugTS-P
MNIST	1885.9 $\pm$ 261.8	1537.8 $\pm$ 99.9	1564.4 $\pm$ 218.2	1451.0 $\pm$ 58.2	<b>1419.1 <math>\pm</math> 45.8</b>
Mushroom	138.0 $\pm$ 31.2	122.5 $\pm$ 15.7	119.1 $\pm$ 26.4	113.0 $\pm$ 21.4	<b>107.4 <math>\pm</math> 19.0</b>
Statlog	204.0 $\pm$ 32.0	180.9 $\pm$ 26.8	179.1 $\pm$ 24.3	163.8 $\pm$ 41.9	<b>150.8 <math>\pm</math> 20.5</b>
Adult	2151.2 $\pm$ 5.6	2150.6 $\pm$ 7.3	2138.8 $\pm$ 10.7	2140.4 $\pm$ 7.7	<b>2136.6 <math>\pm</math> 6.9</b>
Coverttype	2955.0 $\pm$ 37.4	2941.8 $\pm$ 21.3	2947.0 $\pm$ 55.9	2935.9 $\pm$ 47.0	<b>2925.2 <math>\pm</math> 41.2</b>
Magic Telescope	2002.1 $\pm$ 22.1	1995.5 $\pm$ 10.2	1997.9 $\pm$ 17.6	1977.4 $\pm$ 23.2	<b>1969.9 <math>\pm</math> 23.9</b>

Table 6: Total reward averaged over 8 repeated experiments.

	Greedy	MCDropout	NeuralTS	PlugTS	PlugTS-P
Yahoo!R6B	39632.0 $\pm$ 1122.4	47929.1 $\pm$ 476.6	52235.1 $\pm$ 1809.3	52381.8 $\pm$ 540.2	<b>53897.8 <math>\pm</math> 1012.1</b>

Table 7: Overall statistics of classification datasets.

	Mushroom	Statlog	Adult	Coverttype	Magic Telescope
# Instances	8,124	58,000	48,842	581,012	19,020
# Classes	2	7	2	7	2
# Attributes	22	9	14	54	22

**C Details of Experiments**

**C.1 Description of Recommender Simulator** The online recommender system simulator is described in Algorithm 5.

**C.2 Experimental Results** For MCDropout, we set the dropout probability as 0.1 and repeat the forward process three times while predicting. For NeuralTS, PlugTS and PlugTS-P, we use a grid search on  $\lambda \in \{1, 0.1, 0.01\}$  and  $\nu \in \{10^{-3}, 10^{-4}, 10^{-5}\}$ . All experiments are repeated 8 times, and the average and standard error are reported in Table 6 and Table 5. The **Bold Faced** data is the top performance over 8 experiments.

**C.3 Introduction of Classification Datasets** The overall statistics of the classification datasets is shown in Table 7. In the bandit problem setting, the agent obtains reward 1 if the correct class is selected, and 0 otherwise.