
Neural Thompson Sampling

Anonymous Author(s)

Affiliation
Address
email

Abstract

Thompson Sampling (TS) is one of the most effective algorithms for solving contextual multi-armed bandit problems. In this paper, we propose a new algorithm, called Neural Thompson Sampling, which adapts deep neural networks for both exploration and exploitation. At the core of our algorithm is a novel posterior distribution of the reward, where its mean is the neural network approximator, and its variance is built upon the neural tangent features of the corresponding neural network. We prove that, provided the underlying reward function is bounded, the proposed algorithm is guaranteed to achieve a cumulative regret of $\mathcal{O}(T^{1/2})$, which matches the regret of other contextual bandit algorithms in terms of total round number T . Experimental comparisons with other benchmark bandit algorithms on various data sets corroborate our theory.

1 Introduction

The stochastic multi-armed bandit [12, 33] has been extensively studied, as an important model to optimize the trade-off between exploration and exploitation in sequential decision making. Among its many variants, the contextual bandit is widely used in real-world applications such as recommendation [35], advertising [22], robotic control [39], and healthcare [23].

In each round of a contextual bandit, the agent observes a feature vector (the “context”) for each of the K arms, pulls one of them, and in return receives a scalar reward. The goal is to maximize the cumulative reward, or minimize regret (to be defined later), in a total of T rounds. To do so, the agent must find a trade-off between exploration and exploitation. One of the most effective and widely used techniques is *Thompson Sampling*, or TS [47]. The basic idea is to compute the posterior distribution of each arm being optimal for the present context, and sample an arm from this distribution. TS is often easy to implement, and has found great success in practice [14, 22, 28, 45].

Recently, a series of work has applied TS or its variants to explore in contextual bandits with neural network models [11, 31, 38, 43]. Riquelme et al. [43] proposed NeuralLinear, which maintains a neural network and chooses the best arm in each round according to a Bayesian linear regression on top of the last network layer. Kveton et al. [31] proposed DeepFPL, which trains a neural network based on perturbed training data and chooses the best arm in each round based on the neural network output. Similar approaches have also been used in more general reinforcement learning problem [e.g., 10, 20, 37, 41]. Despite the reported empirical success, strong regret guarantees for TS remain limited to relatively simple models, under fairly restrictive assumptions on the reward function. Examples are linear functions [2, 4, 30, 44], generalized linear functions [31, 44], or functions with small RKHS norm induced by a properly selected kernel [15].

In this paper, we provide, to the best of our knowledge, the first near-optimal regret bound for neural network-based Thompson Sampling. Our contributions are threefold. First, we propose a new algorithm, *Neural Thompson Sampling (NeuralTS)*, to incorporate TS exploration with neural networks. It differs from NeuralLinear [43] by considering weight uncertainty in all layers, and from

38 other neural network-based TS implementations [11, 31] by directly sampling the estimated reward
 39 from the posterior (as opposed to sampling parameters).

40 Second, we give a regret analysis for the algorithm, and obtain an $\tilde{\mathcal{O}}(\tilde{d}\sqrt{T})$ regret, where \tilde{d} is the
 41 *effective dimension* and T is the number of rounds. This result is comparable to previous bounds
 42 when specialized to the simpler, linear setting where the effective dimension coincides with the
 43 feature dimension [4, 15].

44 Finally, we corroborate the analysis with an empirical evaluation of the algorithm on several bench-
 45 marks. Experimental results show that NeuralTS yields competitive performance, in comparison with
 46 state-of-the-art baselines, thus suggest its practical value in addition to strong theoretical guarantees.

47 **Notation:** We denote all scalars as lower case letters while constants are denoted as upper case
 48 letters. Vectors are denoted by lower case bold face letters \mathbf{x} and matrices are denoted by upper
 49 case bold face letters \mathbf{A} . we denote by $[k]$ the set $\{1, 2, \dots, k\}$ for all positive integer k . For two
 50 non-negative sequence $\{a_n\}, \{b_n\}$, $a_n = \mathcal{O}(b_n)$ means that there exists a positive constant C such
 51 that $a_n \leq Cb_n$, and we use $\tilde{\mathcal{O}}(\cdot)$ to hide the log factor in $\mathcal{O}(\cdot)$. For norms of vectors and matrices,
 52 we denote by $\|\cdot\|_2$ the Euclidean norm of vectors and the spectral norm of matrices, and denote by
 53 $\|\cdot\|_F$ the Frobenius norm of a matrix.

54 2 Problem settings and algorithm

55 In this work, we consider contextual K -armed bandits, where the total number of rounds T is known.
 56 At round $t \in [T]$, the agent observes K contextual vectors $\{\mathbf{x}_{t,k} \in \mathbb{R}^d \mid k \in [K]\}$. Then the agent
 57 selects an arm a_t and receives a reward r_{t,a_t} . Our goal is to minimize the following pseudo regret:

$$R_T = \mathbb{E} \left[\sum_{t=1}^T (r_{t,a_t^*} - r_{t,a_t}) \right], \quad (2.1)$$

58 where a_t^* is an optimal arm at round t that has the maximum expected reward: $a_t^* = \arg\max_{a \in [K]} \mathbb{E}[r_{t,a}]$. To estimate the unknown reward given a contextual vector \mathbf{x} , we use a fully
 59 connected neural network $f(\mathbf{x}; \boldsymbol{\theta})$ of depth $L \geq 2$, defined recursively by
 60

$$\begin{aligned} f_1 &= \mathbf{W}_1 \mathbf{x}, \\ f_l &= \mathbf{W}_l \text{ReLU}(f_{l-1}), \quad 2 \leq l \leq L, \\ f(\mathbf{x}; \boldsymbol{\theta}) &= \sqrt{m} f_L, \end{aligned} \quad (2.2)$$

61 where $\text{ReLU}(x) := \max\{x, 0\}$ is the widely used Rectified Linear Unit (ReLU) activation function
 62 and $\boldsymbol{\theta} = (\text{vec}(\mathbf{W}_1); \dots; \text{vec}(\mathbf{W}_L)) \in \mathbb{R}^p$ is the collection of parameters of the neural network,
 63 $p = dm + m^2(L-2) + m$.

64 Our proposed Neural Thompson Sampling is given in Algorithm 1. It maintains a Gaussian dis-
 65 tribution for each arm's reward. When selecting an arm, it samples the reward of each arm from
 66 the reward's posterior distribution, and then pulls the greedy arm (lines 4–8). Once the reward is
 67 observed, it updates the posterior (lines 9 & 10). The mean of the posterior distribution is set to the
 68 output of the neural network, whose parameter is the solution to the following minimization problem:

$$\min_{\boldsymbol{\theta}} L(\boldsymbol{\theta}) = \sum_{i=1}^t [f(\mathbf{x}_{i,a_i}; \boldsymbol{\theta}) - r_{i,a_i}]^2 / 2 + m\lambda \|\boldsymbol{\theta} - \boldsymbol{\theta}_0\|_2^2 / 2. \quad (2.3)$$

69 We can see that (2.3) is an ℓ_2 -regularized square loss, where the regularization term centers at the
 70 randomly initialized network parameter $\boldsymbol{\theta}_0$. We adapt gradient descent to solve (2.3) with step size η
 71 and total number of iterations J .

72 A few observations about our algorithm are in place. First, compared to typical ways of implementing
 73 Thompson Sampling with neural networks, NeuralTS samples from the posterior distribution of the
 74 *scalar reward*, instead of the network parameters. It is therefore simpler and more efficient, as the
 75 number of parameters in practice can be large.

76 Second, the algorithm maintains the posterior distributions related to parameters of all layers of the
 77 network, as opposed to the last layer only [43]. This difference is crucial in our regret analysis. It

Algorithm 1 Neural Thompson Sampling (NeuralTS)

Input: Number of rounds T , exploration variance ν , network width m , regularization parameter λ .

- 1: Set $\mathbf{U}_0 = \lambda \mathbf{I}$
- 2: Initialize $\boldsymbol{\theta}_0 = (\text{vec}(\mathbf{W}_1); \dots; \text{vec}(\mathbf{W}_L)) \in \mathbb{R}^p$, where for each $1 \leq l \leq L - 1$, $\mathbf{W}_l = (\mathbf{W}, \mathbf{0}; \mathbf{0}, \mathbf{W})$, each entry of \mathbf{W} is generated independently from $N(0, 4/m)$; $\mathbf{W}_L = (\mathbf{w}^\top, -\mathbf{w}^\top)$, each entry of \mathbf{w} is generated independently from $N(0, 2/m)$.
- 3: **for** $t = 1, \dots, T$ **do**
- 4: **for** $k = 1, \dots, K$ **do**
- 5: $\sigma_{t,k}^2 = \lambda \mathbf{g}^\top(\mathbf{x}_{t,k}; \boldsymbol{\theta}_{t-1}) \mathbf{U}_{t-1}^{-1} \mathbf{g}(\mathbf{x}_{t,k}; \boldsymbol{\theta}_{t-1})/m$
- 6: Sample estimated reward $\tilde{r}_{t,k} \sim \mathcal{N}(f(\mathbf{x}_{t,k}; \boldsymbol{\theta}_{t-1}), \nu^2 \sigma_{t,k}^2)$
- 7: **end for**
- 8: Pull arm a_t and receive reward r_{t,a_t} , where $a_t = \text{argmax}_a \tilde{r}_{t,a}$
- 9: Set $\boldsymbol{\theta}_t$ to be the output of gradient descent for solving (2.3)
- 10: $\mathbf{U}_t = \mathbf{U}_{t-1} + \mathbf{g}(\mathbf{x}_{t,a_t}; \boldsymbol{\theta}_t) \mathbf{g}(\mathbf{x}_{t,a_t}; \boldsymbol{\theta}_t)^\top / m$
- 11: **end for**

78 allows us to build a connection between Algorithm I and recent work about deep learning theory
79 [6, 13], in order to obtain theoretical guarantees as will be shown in the next section.

80 Third, different from linear or kernelized TS [4, 15], whose posterior can be exactly computed in
81 closed forms, NeuralTS solves a non-convex optimization problem (2.3) by gradient descent. This
82 difference requires additional techniques to prove regret bound of the algorithm. We also note that
83 the more commonly used *stochastic* gradient descent can be used to solve the optimization problem
84 with a similar theoretical guarantee [6, 17, 52]. For simplicity of exposition, we will focus on the
85 exact gradient descent approach.

86 **3 Regret analysis**

87 In this section, we provide a regret analysis of NeuralTS. We assume that there exists an unknown
88 reward function h such that for any $1 \leq t \leq T$ and $1 \leq k \leq K$,

$$r_{t,k} = h(\mathbf{x}_{t,k}) + \xi_{t,k}, \quad \text{with } |h(\mathbf{x}_{t,k})| \leq 1$$

89 where $\{\xi_{t,k}\}$ forms an R -sub-Gaussian martingale difference sequence with constant $R > 0$, i.e.,
90 $\mathbb{E}[\exp(\lambda \xi_{t,k}) | \xi_{1:t-1,k}, \mathbf{x}_{1:t,k}] \leq \exp(\lambda^2 R^2)$ for all $\lambda \in \mathbb{R}$. Such an assumption on the noise sequence
91 is widely adapted in contextual bandit literature [4, 12, 15, 16, 33, 48]. Note that we make *no*
92 assumption on h other than boundedness.

93 Next, we provide necessary background on the neural tangent kernel (NTK) theory [26], which plays
94 a crucial role in our analysis. In the remaining of the analysis, we denote by $\{\mathbf{x}^i\}_{i=1}^{TK}$ the set of
95 observed contexts of all arms and all rounds: $\{\mathbf{x}_{t,k}\}_{1 \leq t \leq T, 1 \leq k \leq K}$ where $i = K(t-1) + k$.

96 **Definition 3.1** (Jacot et al. [26]). Define

$$\begin{aligned} \widetilde{\mathbf{H}}_{i,j}^{(1)} &= \boldsymbol{\Sigma}_{i,j}^{(1)} = \langle \mathbf{x}^i, \mathbf{x}^j \rangle, \quad \mathbf{A}_{i,j}^{(l)} = \begin{pmatrix} \boldsymbol{\Sigma}_{i,i}^{(l)} & \boldsymbol{\Sigma}_{i,j}^{(l)} \\ \boldsymbol{\Sigma}_{i,j}^{(l)} & \boldsymbol{\Sigma}_{j,j}^{(l)} \end{pmatrix}, \\ \boldsymbol{\Sigma}_{i,j}^{(l+1)} &= 2\mathbb{E}_{(u,v) \sim N(\mathbf{0}, \mathbf{A}_{i,j}^{(l)})} \max\{u, 0\} \max\{v, 0\}, \\ \widetilde{\mathbf{H}}_{i,j}^{(l+1)} &= 2\widetilde{\mathbf{H}}_{i,j}^{(l)} \mathbb{E}_{(u,v) \sim N(\mathbf{0}, \mathbf{A}_{i,j}^{(l)})} \mathbb{1}(u \geq 0) \mathbb{1}(v \geq 0) + \boldsymbol{\Sigma}_{i,j}^{(l+1)}. \end{aligned}$$

97 Then, $\mathbf{H} = (\widetilde{\mathbf{H}}^{(L)} + \boldsymbol{\Sigma}^{(L)})/2$ is called the neural tangent kernel matrix on the context set.

98 The NTK technique builds a connection between deep neural networks and kernel methods. It enables
99 us to adapt some complexity measures for kernel methods to describe the complexity of the neural
100 network, as given by the following definition.

101 **Definition 3.2.** The *effective dimension* \tilde{d} of matrix \mathbf{H} with regularization parameter λ is defined as

$$\tilde{d} = \frac{\log \det(\mathbf{I} + \mathbf{H}/\lambda)}{\log(1 + TK/\lambda)}.$$

102 **Remark 3.3.** The effective dimension is a metric to describe the actual underlying dimension in the
 103 set of observed contexts, and has been used by Valko et al. [48] for the analysis of kernel UCB. Our
 104 definition here is adapted from Yang and Wang [49], who also considers UCB-based exploration.
 105 Compared with the maximum information gain γ_t used in Chowdhury and Gopalan [15], one can
 106 verify that their Lemma 3 shows that $\gamma_t \geq \log \det(\mathbf{I} + \mathbf{H}/\lambda)/2$. Therefore, γ_t and \tilde{d} are of the same
 107 order up to a ratio of $1/(2 \log(1 + TK/\lambda))$.

108 We will make a regularity assumption on the contexts and the corresponding NTK matrix \mathbf{H} .

109 **Assumption 3.4.** Let \mathbf{H} be defined in Definition 3.1. There exists $\lambda_0 > 0$, such that $\mathbf{H} \succeq \lambda_0 \mathbf{I}$. In
 110 addition, for any $t \in [T]$, $k \in [K]$, $\|\mathbf{x}_{t,k}\|_2 = 1$ and $[\mathbf{x}_{t,k}]_j = [\mathbf{x}_{t,k}]_{j+d/2}$.

111 The assumption that the NTK matrix is positive definite has been considered in prior work on NTK
 112 [7, 17]. The assumption on context $\mathbf{x}_{t,a}$ ensures that the initial output of neural network $f(\mathbf{x}; \theta_0)$ is 0
 113 with the random initialization suggested in Algorithm 1. The condition on \mathbf{x} is easy to satisfy, since
 114 for any context \mathbf{x} , one can always construct a new context $\tilde{\mathbf{x}}$ as $[\mathbf{x}/(\sqrt{2}\|\mathbf{x}\|_2), \mathbf{x}/(\sqrt{2}\|\mathbf{x}\|_2)]^\top$.

115 We are now ready to present the main result of the paper:

116 **Theorem 3.5.** Under Assumption 3.4, set the parameters in Algorithm 1 as $\lambda = 1 + 1/T$, $\nu =$
 117 $B + R\sqrt{\tilde{d}\log(1 + TK/\lambda) + 2 + 2\log(1/\delta)}$ where $B = \max\left\{1/(22e\sqrt{\pi}), \sqrt{2\mathbf{h}^\top \mathbf{H}^{-1} \mathbf{h}}\right\}$ with
 118 $\mathbf{h} = (h(\mathbf{x}^1), \dots, h(\mathbf{x}^{TK}))^\top$, and R is the sub-Gaussian parameter. For gradient descent in line
 119 9 of Algorithm 1, set the step size as $\eta = C_1(m\lambda + mL)^{-1}$ and the number of iterations as
 120 $J = (1 + LT/\lambda)(C_2 + \log(T^3L\lambda^{-1}\log(1/\delta)))/C_1$ for some positive constant C_1, C_2 . If the
 121 network width m satisfies:

$$m \geq \text{poly}\left(\lambda, T, K, L, \log(1/\delta), \lambda_0^{-1}\right),$$

122 then, with probability at least $1 - \delta$, the regret of Algorithm 1 is bounded as

$$R_T \leq C_2(1 + c_T)\nu\sqrt{2\lambda L(\tilde{d}\log(1 + TK) + 1)T} + (4 + C_3(1 + c_T)\nu L)\sqrt{2\log(3/\delta)T} + 5,$$

123 where C_2, C_3 are absolute constants, and $c_T = \sqrt{4\log T + 2\log K}$.

124 **Remark 3.6.** The definition B in Theorem 3.5 is inspired by the RKHS norm of the reward function
 125 defined in Chowdhury and Gopalan [15]. It can be verified that B is an absolute constant as long as
 126 the reward function h belongs to the function space induced by NTK.

127 **Remark 3.7.** Theorem 3.5 implies the regret of NeuralTS is on the order of $\tilde{O}(\tilde{d}T^{1/2})$. This result
 128 matches the state-of-the-art regret bound in Chowdhury and Gopalan [15], Agrawal and Goyal
 129 [4], Zhou et al. [51], Kveton et al. [31].

130 4 Proof of the main theorem

131 This section sketches the proof of Theorem 3.5, with supporting lemmas and technical details provided
 132 in Appendix B. While the proof roadmap is similar to previous work on Thompson Sampling [e.g.,
 133 4, 15, 29, 31], our proof needs to carefully track the approximation error of neural networks for
 134 approximating the reward function. To control the approximation error, the following condition on
 135 the neural network width is required in many technical lemmas.

136 **Condition 4.1.** For the network width m and step size η , there exists a set of positive constants
 137 $\{C_{m,1}, C_{m,2}, \dots, C_{m,5}\}$ such that

$$\begin{aligned} m^2 &\geq C_{m,1}L^{-3/2}\lambda^{1/2}[\log(TKL^2/\delta)]^{3/2}, \\ 2\sqrt{T/\lambda} &\leq C_{m,2}\min\left\{mL^{-6}[\log m]^{-3/2}, (m^2(\lambda\eta)^2L^{-6}T^{-1}(\log m)^{-1})^{3/8}\right\}, \\ m^{1/6} &\geq C_{m,3}\sqrt{\log m}L^{7/2}T^{7/6}\lambda^{-7/6}(1 + \sqrt{T/\lambda}) \\ m &\geq C_{m,4}T^6K^6L^6\log(TKL/\delta)\max\{\lambda_0^{-4}, 1\} \\ \eta &= C_{m,5}(m\lambda + mL)^{-1}. \end{aligned}$$

138 We define \mathcal{F}_t as the σ -algebra $\mathcal{F}_t = \sigma(\boldsymbol{\theta}_0, \tilde{r}_{1,1}, \dots, \tilde{r}_{t-1,K}, r_{1,a_1}, \dots, r_{t-1,a_{t-1}})$. For any t , we
 139 define an event \mathcal{E}_t^σ as follows

$$\mathcal{E}_t^\sigma = \{\omega \in \mathcal{F}_{t+1} : \forall k \in [K], \quad |\tilde{r}_{t,k} - f(\mathbf{x}_{t,k}; \boldsymbol{\theta}_{t-1})| \leq c_t \nu \sigma_{t,k}\}, \quad (4.1)$$

140 where $c_t = \sqrt{4 \log t + 2 \log K}$. Under event \mathcal{E}_t^σ , the difference between the sampled reward $\tilde{r}_{t,k}$ and
 141 the estimated mean reward $f(\mathbf{x}_{t,k}; \boldsymbol{\theta}_{t-1})$ can be controlled by the reward's posterior variance.

142 We also define an event \mathcal{E}_t^μ as follows

$$\mathcal{E}_t^\mu = \{\omega \in \mathcal{F}_t : \forall k \in [K], \quad |f(\mathbf{x}_{t,k}; \boldsymbol{\theta}_{t-1}) - h(\mathbf{x}_{t,k})| \leq \nu \sigma_{t,k} + \epsilon(m)\}, \quad (4.2)$$

143 where $\epsilon(m)$ is defined as

$$\begin{aligned} \epsilon(m) &= \epsilon_p(m) + C_{\epsilon,1}(1 - \eta m \lambda)^J \sqrt{TL/\lambda} \\ \epsilon_p(m) &= C_{\epsilon,2} T^{2/3} m^{-1/6} \lambda^{-2/3} L^3 \sqrt{\log m} + C_{\epsilon,3} m^{-1/6} \sqrt{\log m} L^4 T^{5/3} \lambda^{-5/3} (1 + \sqrt{T/\lambda}) \\ &\quad + C_{\epsilon,4} \left(B + R \sqrt{\log \det(\mathbf{I} + \mathbf{H}/\lambda) + 2 + 2 \log(1/\delta)} \right) \sqrt{\log m} T^{7/6} m^{-1/6} \lambda^{-2/3} L^{9/2}, \end{aligned} \quad (4.3)$$

144 and $\{C_{\epsilon,i}\}_{i=1}^4$ are absolute constants. Under event \mathcal{E}_t^μ , the estimated mean reward $f(\mathbf{x}_{t,k}; \boldsymbol{\theta}_{t-1})$
 145 based on the neural network is similar to the true expected reward $h(\mathbf{x}_{t,k})$. Note that the additional
 146 term $\epsilon(m)$ is the approximate error of the neural networks for approximating the true reward function.
 147 This is a key difference in our proof from previous regret analysis of Thompson Sampling [4, 15],
 148 where there is no approximation error.

149 The following two lemmas show that both events \mathcal{E}_t^σ and \mathcal{E}_t^μ happen with high probability.

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

151 **Lemma 4.3.** Under Condition 4.1, $\Pr(\forall t \in [T], \mathcal{E}_t^\mu) \geq 1 - \delta$.

152 The next lemma gives a lower bound of the probability that the sampled reward \tilde{r} is larger than true
 153 reward up to the approximation error $\epsilon(m)$.

154 **Lemma 4.4.** For any $t \in [T]$ and for all $k \in [K]$, we have

$$\Pr(\tilde{r}_{t,k} + \epsilon(m) > h(\mathbf{x}_{t,k}) | \mathcal{F}_t, \mathcal{E}_t^\mu) \geq \frac{1}{4e\sqrt{\pi}}.$$

155 Following Agrawal and Goyal [4], for any time t , we divide the arms into two groups: saturated and
 156 unsaturated arms, based on whether the standard deviation of the estimates for an arm is smaller than
 157 the standard deviation for the optimal arm or not. Note that the optimal arm is included in the group
 158 of unsaturated arms. More specifically, we define the set of saturated arms S_t as follows

$$S_t = \{k | k \in [K], h(\mathbf{x}_{t,a_t^*}) - h(\mathbf{x}_{t,k}) \geq (1 + c_t) \nu \sigma_{t,k} + 2\epsilon(m)\}. \quad (4.4)$$

159 Note that here we also take the approximate error $\epsilon(m)$ into consideration when we define the
 160 saturated arms, which differs from the definition of saturated arms in classical Thompson Sampling
 161 literature [4, 15]. It is now easy to show that the immediate regret of playing an unsaturated arm can
 162 be bounded by the standard deviation plus the approximation error $\epsilon(m)$.

163 The following lemma shows that the probability of pulling a saturated arm is small in Algorithm 1.

164 **Lemma 4.5.** Let a_t be the arm pulled at round $t \in [T]$. Then,

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

165 The next lemma bounds the expectation of the regret at each round conditioned on \mathcal{E}_t^μ .

166 **Lemma 4.6.** Under Condition 4.1, with probability at least $1 - \delta$, we have for all $t \in [T]$ that

$$\mathbb{E}[h(\mathbf{x}_{t,a_t^*}) - h(\mathbf{x}_{t,a_t}) | \mathcal{F}_t, \mathcal{E}_t^\mu] \leq C(1 + c_t) \nu \sqrt{L} \mathbb{E}[\min\{\sigma_{t,a_t}, 1\} | \mathcal{F}_t, \mathcal{E}_t^\mu] + 4\epsilon(m) + 2t^{-2},$$

167 where C is an absolute constant.

168 Based on Lemma 4.6 we define

$$\begin{aligned}\bar{\Delta}_t &:= (h(\mathbf{x}_{t,a_t^*}) - h(\mathbf{x}_{t,a_t})) \mathbb{1}(\mathcal{E}_t^\mu) \\ X_t &:= \bar{\Delta}_t - (C_\Delta(1 + c_t)\nu\sqrt{L} \min\{\sigma_{t,a_t}, 1\} + 4\epsilon(m) + 2t^{-2}), \quad Y_t = \sum_{i=1}^t X_i,\end{aligned}\quad (4.5)$$

169 where C_Δ is the same with constant C in Lemma 4.6. By Lemma 4.6, we can verify that with
170 probability at least $1 - \delta$, $\{Y_t\}$ forms a super martingale sequence since $\mathbb{E}(Y_t - Y_{t-1}) = \mathbb{E}X_t \leq 0$.
171 By Azuma-Hoeffding inequality [24], we can prove the following lemma.

172 **Lemma 4.7.** Under Condition 4.1, we have, with probability at least $1 - \delta$, that

$$\begin{aligned}\sum_{i=1}^T \bar{\Delta}_i &\leq 4T\epsilon(m) + \pi^2/3 + C_1(1 + c_T)\nu\sqrt{L} \sum_{i=1}^T \min\{\sigma_{t,a_t}, 1\} \\ &\quad + (4 + C_2(1 + c_T)\nu L + 4\epsilon(m))\sqrt{2\log(1/\delta)T},\end{aligned}$$

173 where C_1, C_2 are absolute constants.

174 The last lemma is used to control $\sum_{i=1}^T \min\{\sigma_{t,a_t}, 1\}$ in Lemma 4.7.

175 **Lemma 4.8.** Under Condition 4.1, with probability at least $1 - \delta$, it holds that

$$\sum_{i=1}^T \min\{\sigma_{t,a_t}, 1\} \leq \sqrt{2\lambda T(\tilde{d}\log(1 + TK) + 1)} + CT^{13/6}\sqrt{\log mm^{-1/6}}\lambda^{-2/3}L^{9/2},$$

176 where C is an absolute constant.

177 With the above lemmas, we are ready to prove Theorem 3.5

178 *Proof of Theorem 3.5.* Under Condition 4.1, by Lemma 4.3, \mathcal{E}_t^μ holds for all $t \in [T]$ with probability
179 at least $1 - \delta$. Therefore, with probability at least $1 - \delta$, we have

$$\begin{aligned}R_T &= \sum_{i=1}^T (h(\mathbf{x}_{t,a_t^*}) - h(\mathbf{x}_{t,a_t})) \mathbb{1}(\mathcal{E}_t^\mu) \\ &\leq 4T\epsilon(m) + \frac{\pi^2}{3} + \bar{C}_1(1 + c_T)\nu\sqrt{L} \sum_{i=1}^T \min\{\sigma_{t,a_t}, 1\} \\ &\quad + (4 + \bar{C}_2(1 + c_T)\nu L + 4\epsilon(m))\sqrt{2\log(1/\delta)T} \\ &\leq \bar{C}_1(1 + c_T)\nu\sqrt{L} \left(\sqrt{2\lambda T(\tilde{d}\log(1 + TK) + 1)} + \bar{C}_3 T^{13/6} \sqrt{\log mm^{-1/6}} \lambda^{-2/3} L^{9/2} \right) \\ &\quad + \frac{\pi^2}{3} + 4T\epsilon(m) + 4\epsilon(m)\sqrt{2\log(1/\delta)T} + (4 + \bar{C}_2(1 + c_T)\nu L)\sqrt{2\log(1/\delta)T}, \\ &= \bar{C}_1(1 + c_T)\nu\sqrt{L} \left(\sqrt{2\lambda T(\tilde{d}\log(1 + TK) + 1)} + \bar{C}_3 T^{13/6} \sqrt{\log mm^{-1/6}} \lambda^{-2/3} L^{9/2} \right) \\ &\quad + \frac{\pi^2}{3} + \epsilon_p(m)(4T + \sqrt{2\log(1/\delta)T}) + (4 + \bar{C}_2(1 + c_T)\nu L)\sqrt{2\log(1/\delta)T} \\ &\quad + C_{\epsilon,1}(1 - \eta m\lambda)^J \sqrt{TL/\lambda}(4T + \sqrt{2\log(1/\delta)T}),\end{aligned}$$

180 where $\bar{C}_1, \bar{C}_2, \bar{C}_3 > 0$ are absolute constants, the first inequality is due to Lemma 4.7 and the second
181 inequality is due to Lemma 4.8. The third equation is from (4.3). By setting $\eta = C_4(m\lambda + mL\lambda)^{-1}$
182 and $J = (1 + LT/\lambda)(\log(24C_{\epsilon,1}) + \log(T^3L\lambda^{-1}\log(1/\delta)))/C_4$, we have

$$C_{\epsilon,1}(1 - \eta m\lambda)^J \sqrt{TL/\lambda}(4T + \sqrt{2\log(1/\delta)T}) \leq \frac{1}{3},$$

183 Then choosing m such that

$$\bar{C}_1 \bar{C}_3 (1 + c_T) \nu T^{13/6} \sqrt{\log m} m^{-1/6} \lambda^{-2/3} L^5 \leq \frac{1}{3}, \quad \epsilon_p(m)(4T + \sqrt{2 \log(1/\delta)T}) \leq \frac{1}{3}.$$

184 R_T can be further bounded by

$$R_T \leq \bar{C}_1 (1 + c_T) \nu \sqrt{2\lambda L(\tilde{d} \log(1 + TK) + 1)T} + (4 + \bar{C}_2 (1 + c_T) \nu L) \sqrt{2 \log(1/\delta)T} + 5.$$

185 Taking union bound over Lemmas 4.3, 4.7 and 4.8, the above inequality holds with probability $1 - 3\delta$.
186 By replacing δ with $\delta/3$ and rearranging terms, we complete the proof. \square

187 5 Experiments

188 This section gives an empirical evaluation of our algorithm in several public benchmark datasets,
189 including adult, covtype, magic telescope, mushroom and shuttle, all from UCI [18], as
190 well as MNIST [34]. The algorithm is compared to several typical baselines: linear and kernelized
191 Thompson Sampling [4, 15], linear and kernelized UCB [16, 48], BootstrapNN [42, 43], and ϵ -Greedy
192 for neural networks. BootstrapNN trains multiple neural networks with subsampled data, and at each
193 step pulls the greedy action based on a randomly selected network. It has been proposed as a way to
194 approximate Thompson Sampling [40, 42].

195 5.1 Experiment setup

196 To transform these classification problems into multi-armed bandit problems, we adapt the disjoint
197 models [35] to build a context feature vector for each arm. In detail, for a classification problem
198 with $\mathbf{x} \in \mathbb{R}^d$ as the input feature vector and k target classes, we build the context feature vector
199 with dimension kd as $\mathbf{x}_1 = (\mathbf{x}; \mathbf{0}; \dots; \mathbf{0}), \mathbf{x}_2 = (\mathbf{0}; \mathbf{x}; \dots; \mathbf{0}), \dots, \mathbf{x}_k = (\mathbf{0}; \mathbf{0}; \dots; \mathbf{x})$. Then the
200 algorithm generates a set of predicted reward following Algorithm 1 and pulls the greedy arm. For
201 these classification problems, if the algorithm selects a correct class by pulling the corresponding arm,
202 it will receive a reward as 1, otherwise 0. The cumulative regret over time horizon T is measured by
203 the total mistakes made by the algorithm.

204 We set the time horizon of our algorithm to 10 000 for all data sets, except for mushroom which
205 contains only 8 124 data. In order to speed up training, instead of calculating the inverse matrix
206 of \mathbf{U} , we use the inverse of the diagonal elements of \mathbf{U} as an approximation of \mathbf{U}^{-1} . Also, since
207 calculating the kernel matrix is very expensive, we stop training at $t = 1000$ and keep evaluating the
208 performance in the rest of the time, similar to previous work [43, 51].

209 In the experiments, we shuffle all datasets randomly, and normalize the features so that their ℓ_2 -norm
210 is unity. One-hidden-layer neural networks with 100 neurons are used. During posterior updating,
211 gradient descent is run for 100 iterations with learning rate 0.001. For BootstrapNN, we use 10
212 identical networks, and to train each network, data point at each round has probability 0.8 to be
213 included for training ($p = 10, q = 0.8$ in the original paper [46]). For ϵ -Greedy, we tune ϵ with a
214 grid search on $\{0.01, 0.05, 0.1\}$. For (λ, ν) used in linear and kernel UCB / Thompson Sampling, we
215 set $\lambda = 1$ following previous works [4, 15], and do a grid search of $\nu \in \{1, 0.1, 0.01\}$ to select the
216 parameter with best performance. For the Neural UCB / Thompson Sampling methods, we use a grid
217 search on $\lambda \in \{1, 10^{-1}, 10^{-2}, 10^{-3}\}$ and $\nu \in \{10^{-1}, 10^{-2}, 10^{-3}, 10^{-4}, 10^{-5}\}$. All experiments
218 are repeated 8 times, and the average and standard error are reported.

219 5.2 Results

220 The experiment results are shown in Figure 1, and a few observations are in place. First, Neural
221 Thompson Sampling's performance is among the best in 6 datasets, and is significantly better than all
222 other baselines in 2 of them. Second, the function class used by an algorithm is important. Those
223 with linear representations tend to perform worse due to the nonlinearity of rewards in the data. Third,
224 Thompson Sampling is competitive with, and sometimes better than, other exploration strategies
225 with the same function class, in particular when neural networks are used. More detailed results are
226 provided in Appendix A.

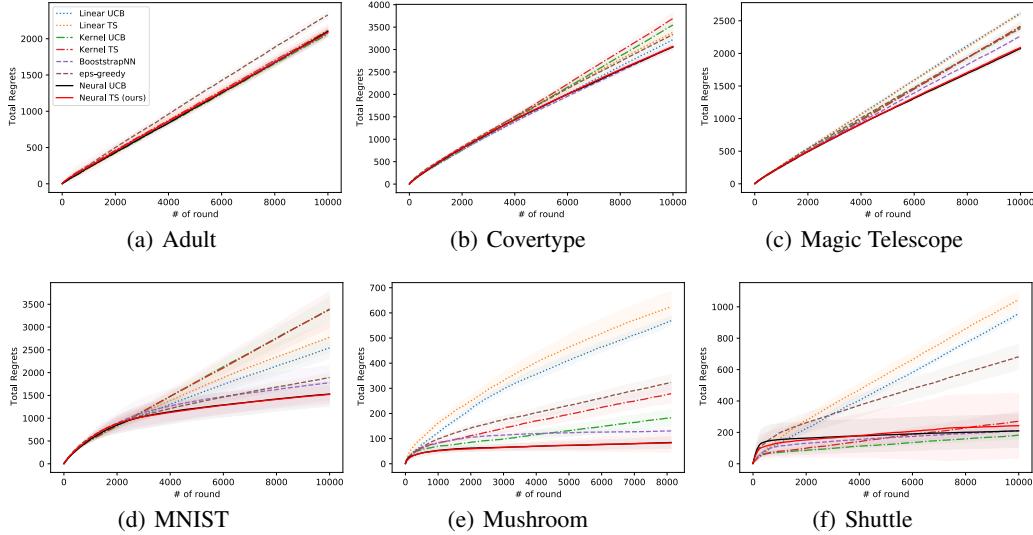


Figure 1: Comparison of Neural Thompson Sampling and baselines on UCI dataset and MNIST dataset. The total regret measures cumulative classification errors made by an algorithm. Results are averaged over multiple runs with standard errors shown as shaded areas.

227 6 Related work

228 Thompson Sampling was proposed as an exploration heuristic almost nine decades ago [47], and
229 has received significant interest in the last decade. Previous works related to the present paper are
230 discussed in the introduction, and are not repeated here.

231 Upper confidence bound or UCB [3, 9, 32] is a widely used alternative to Thompson Sampling for
232 exploration. This strategy is shown to achieve near-optimal regrets in a range of settings, such as linear
233 bandits [1, 8, 16], generalized linear bandits [19, 27, 36], and kernelized contextual bandits [48].

234 Neural networks are increasingly used in contextual bandits. In addition to those mentioned earlier [11,
235 31, 38, 43], Zahavy and Mannor [50] used a deep neural network to provide a feature mapping and
236 explored only at the last layer. Schwenk and Bengio [46] proposed an algorithm by boosting the
237 estimation of multiple deep neural networks. While these methods all show promise empirically,
238 no regret guarantees are known. Recently, Zhou et al. [51] proposed a neural UCB algorithm with
239 near-optimal regret; Foster and Rakhlin [21] proposed a special regression oracle for contextual
240 bandits with a general function class (including neural networks) along with theoretical analysis.
241 Both these works are based on UCB exploration, while this paper focuses on Thompson Sampling.

242 7 Conclusion

243 In this paper, we adapt Thompson Sampling to neural networks. Building on recent advances in
244 deep neural networks theory, we are able to show that the proposed algorithm, NeuralTS, enjoys a
245 $\tilde{\mathcal{O}}(\tilde{d}T^{1/2})$ regret bound. We also show the algorithm works well empirically on benchmark problems,
246 in comparison with multiple strong baselines.

247 The promising results suggest a few interesting directions for future research. First, our analysis needs
248 NeuralTS to perform multiple gradient descent steps to train the neural network in each round. It is
249 interesting to analyze the case where NeuralTS only performs one gradient descent step in each round,
250 and in particular, the trade-off between optimization precision and regret minimization. Second, when
251 the number of arms is finite, $\tilde{\mathcal{O}}(\sqrt{dT})$ regret has been established for parametric bandits with linear
252 and generalized linear reward functions. It is an open problem how to adapt NeuralTS to achieve
253 the same rate. Third, recent work [5] suggested that neural networks may behave differently from a
254 neural tangent kernel under some parameter regimes. It is interesting to investigate whether similar
255 results hold for neural contextual bandit algorithms like NeuralTS.

256 **Broader Impact**

257 The majority of the paper is theoretical, without immediate intended applications to any domain. That
258 said, balancing exploitation and exploration is of fundamental importance to an Artificial Intelligence
259 system, since it reflects how the AI agent efficiently acquires knowledge from interactions with the
260 environment. A more efficient exploration strategy has the potential to reduce unnecessary risk in this
261 process. On the other hand, exploration comes at the unavoidable cost of choosing suboptimal actions
262 in some rounds, which should be treated with care, as it can lead to ethical or fairness concerns, for
263 example, in healthcare applications.

264 **References**

- 265 [1] Yasin Abbasi-Yadkori, Dávid Pál, and Csaba Szepesvári. Improved algorithms for linear
266 stochastic bandits. In *Advances in Neural Information Processing Systems*, pages 2312–2320,
267 2011.
- 268 [2] Marc Abeille and Alessandro Lazaric. Linear Thompson sampling revisited. In *Proceedings of
269 the 20th International Conference on Artificial Intelligence and Statistics*, pages 176–184, 2017.
- 270 [3] Rajeev Agrawal. Sample mean based index policies by $o(\log n)$ regret for the multi-armed
271 bandit problem. *Advances in Applied Probability*, 27(4):1054–1078, 1995.
- 272 [4] Shipra Agrawal and Navin Goyal. Thompson sampling for contextual bandits with linear
273 payoffs. In *International Conference on Machine Learning*, pages 127–135, 2013.
- 274 [5] Zeyuan Allen-Zhu and Yuanzhi Li. What can resnet learn efficiently, going beyond kernels? In
275 *Advances in Neural Information Processing Systems*, pages 9015–9025, 2019.
- 276 [6] Zeyuan Allen-Zhu, Yuanzhi Li, and Zhao Song. A convergence theory for deep learning via
277 over-parameterization. *arXiv preprint arXiv:1811.03962*, 2018.
- 278 [7] Sanjeev Arora, Simon S Du, Wei Hu, Zhiyuan Li, and Ruosong Wang. Fine-grained analysis
279 of optimization and generalization for overparameterized two-layer neural networks. *arXiv
280 preprint arXiv:1901.08584*, 2019.
- 281 [8] Peter Auer. Using confidence bounds for exploitation-exploration trade-offs. *Journal of Machine
282 Learning Research*, 3(Nov):397–422, 2002.
- 283 [9] Peter Auer, Nicolò Cesa-Bianchi, and Paul Fischer. Finite-time analysis of the multiarmed
284 bandit problem. *Machine Learning*, 47(2–3):235–256, 2002.
- 285 [10] Kamyar Azizzadenesheli, Emma Brunskill, and Animashree Anandkumar. Efficient exploration
286 through Bayesian deep Q-networks. In *Proceedings of the 2018 Information Theory and
287 Applications Workshop*, pages 1–9, 2018.
- 288 [11] Charles Blundell, Julien Cornebise, Koray Kavukcuoglu, and Daan Wierstra. Weight uncertainty
289 in neural network. In *Proceedings of the 32nd International Conference on Machine Learning*,
290 pages 1613–1622, 2015.
- 291 [12] Sébastien Bubeck and Nicolò Cesa-Bianchi. Regret analysis of stochastic and nonstochastic
292 multi-armed bandit problems. *Foundations and Trends in Machine Learning*, 5(1):1–122, 2012.
- 293 [13] Yuan Cao and Quanquan Gu. Generalization bounds of stochastic gradient descent for wide
294 and deep neural networks. In *Advances in Neural Information Processing Systems*, pages
295 10835–10845, 2019.
- 296 [14] Olivier Chapelle and Lihong Li. An empirical evaluation of thompson sampling. In *Advances
297 in neural information processing systems*, pages 2249–2257, 2011.
- 298 [15] Sayak Ray Chowdhury and Aditya Gopalan. On kernelized multi-armed bandits. In *Proceedings
299 of the 34th International Conference on Machine Learning*, pages 844–853, 2017.

- 300 [16] Wei Chu, Lihong Li, Lev Reyzin, and Robert Schapire. Contextual bandits with linear payoff
 301 functions. In *Proceedings of the 14th International Conference on Artificial Intelligence and*
 302 *Statistics*, pages 208–214, 2011.
- 303 [17] Simon S Du, Xiyu Zhai, Barnabas Poczos, and Aarti Singh. Gradient descent provably optimizes
 304 over-parameterized neural networks. *arXiv preprint arXiv:1810.02054*, 2018.
- 305 [18] Dheeru Dua and Casey Graff. UCI machine learning repository, 2017. URL <http://archive.ics.uci.edu/ml>
- 306
- 307 [19] Sarah Filippi, Olivier Cappe, Aurélien Garivier, and Csaba Szepesvári. Parametric bandits: The
 308 generalized linear case. In *Advances in Neural Information Processing Systems*, pages 586–594,
 309 2010.
- 310 [20] Meire Fortunato, Mohammad Gheshlaghi Azar, Bilal Piot, Jacob Menick, Matteo Hessel, Ian
 311 Osband, Alex Graves, Volodymyr Mnih, Rémi Munos, Demis Hassabis, Olivier Pietquin,
 312 Charles Blundell, and Shane Legg. Noisy networks for exploration. In *Proceedings of the 6th*
 313 *International Conference on Learning Representations*, 2018.
- 314 [21] Dylan J Foster and Alexander Rakhlin. Beyond ucb: Optimal and efficient contextual bandits
 315 with regression oracles. *arXiv preprint arXiv:2002.04926*, 2020.
- 316 [22] Thore Graepel, Joaquin Quinonero Candela, Thomas Borchert, and Ralf Herbrich. Web-scale
 317 Bayesian click-through rate prediction for sponsored search advertising in Microsoft’s Bing
 318 search engine. In *Proceedings of the 27th International Conference on Machine Learning*,
 319 pages 13–20, 2010.
- 320 [23] Kristjan Greenewald, Ambuj Tewari, Susan Murphy, and Predag Klasnja. Action centered
 321 contextual bandits. In *Advances in Neural Information Processing Systems 30*, pages 5977–5985,
 322 2017.
- 323 [24] Wassily Hoeffding. Probability inequalities for sums of bounded random variables. *Journal of*
 324 *the American Statistical Association*, 58(301):13–30, 1963.
- 325 [25] Matthew W Hoffman, Bobak Shahriari, and Nando de Freitas. Exploiting correlation and budget
 326 constraints in bayesian multi-armed bandit optimization. *arXiv preprint arXiv:1303.6746*, 2013.
- 327 [26] Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and
 328 generalization in neural networks. In *Advances in Neural Information Processing Systems*,
 329 pages 8571–8580, 2018.
- 330 [27] Kwang-Sung Jun, Aniruddha Bhargava, Robert D. Nowak, and Rebecca Willett. Scalable
 331 generalized linear bandits: Online computation and hashing. In *Advances in Neural Information*
 332 *Processing Systems 30*, pages 99–109, 2017.
- 333 [28] Jaya Kawale, Hung Hai Bui, Branislav Kveton, Long Tran-Thanh, and Sanjay Chawla. Efficient
 334 Thompson sampling for online matrix-factorization recommendation. In *Advances in Neural*
 335 *Information Processing Systems 28*, pages 1297–1305, 2015.
- 336 [29] Tomáš Kocák, Michal Valko, Rémi Munos, and Shipra Agrawal. Spectral Thompson sampling.
 337 In *28th AAAI Conference on Artificial Intelligence*, 2014.
- 338 [30] Tomáš Kocák, Michal Valko, Rémi Munos, and Shipra Agrawal. Spectral Thompson sampling.
 339 In *Proceedings of the 28th AAAI Conference on Artificial Intelligence*, pages 1911–1917, 2014.
- 340 [31] Branislav Kveton, Manzil Zaheer, Csaba Szepesvari, Lihong Li, Mohammad Ghavamzadeh,
 341 and Craig Boutilier. Randomized exploration in generalized linear bandits. In *Proceedings of*
 342 *the 22nd International Conference on Artificial Intelligence and Statistics*, 2020.
- 343 [32] Tze Leung Lai and Herbert Robbins. Asymptotically efficient adaptive allocation rules. *Ad-*
 344 *vances in Applied Mathematics*, 6(1):4–22, 1985.
- 345 [33] Tor Lattimore and Csaba Szepesvári. *Bandit Algorithms*. Cambridge University Press, 2020.

- 346 [34] Yann LeCun, Corinna Cortes, and CJ Burges. Mnist handwritten digit database. *ATT Labs*
 347 *[Online]. Available: <http://yann.lecun.com/exdb/mnist>*, 2, 2010.
- 348 [35] Lihong Li, Wei Chu, John Langford, and Robert E Schapire. A contextual-bandit approach to
 349 personalized news article recommendation. In *Proceedings of the 19th International Conference*
 350 *on World Wide Web*, pages 661–670, 2010.
- 351 [36] Lihong Li, Yu Lu, and Dengyong Zhou. Provably optimal algorithms for generalized linear
 352 contextual bandits. In *Proceedings of the 34th International Conference on Machine Learning-
 353 Volume 70*, pages 2071–2080. JMLR.org, 2017.
- 354 [37] Zachary C. Lipton, Jianfeng Gao, Lihong Li, Xiujun Li, Faisal Ahmed, and Li Deng. BBQ-
 355 networks: Efficient exploration in deep reinforcement learning for task-oriented dialogue
 356 systems. In *Proceedings of the 32nd AAAI Conference on Artificial Intelligence*, pages 5237–
 357 5244, 2018.
- 358 [38] Xiuyuan Lu and Benjamin Van Roy. Ensemble sampling. In *Advances in Neural Information
 359 Processing Systems 30*, pages 3258–3266, 2017.
- 360 [39] Jeffrey Mahler, Florian T Pokorny, Brian Hou, Melrose Roderick, Michael Laskey, Mathieu
 361 Aubry, Kai Kohlhoff, Torsten Kröger, James Kuffner, and Ken Goldberg. Dex-net 1.0: A
 362 cloud-based network of 3d objects for robust grasp planning using a multi-armed bandit model
 363 with correlated rewards. In *2016 IEEE international conference on robotics and automation
 (ICRA)*, pages 1957–1964. IEEE, 2016.
- 365 [40] Ian Osband and Benjamin Van Roy. Bootstrapped thompson sampling and deep exploration.
 366 *arXiv preprint arXiv:1507.00300*, 2015.
- 367 [41] Ian Osband, Charles Blundell, Alexander Pritzel, and Benjamin Van Roy. Deep exploration
 368 via bootstrapped DQN. In *Advances in Neural Information Processing Systems 29*, pages
 369 4026–4034, 2016.
- 370 [42] Ian Osband, Charles Blundell, Alexander Pritzel, and Benjamin Van Roy. Deep exploration via
 371 bootstrapped dqn. In *Advances in neural information processing systems*, pages 4026–4034,
 372 2016.
- 373 [43] Carlos Riquelme, George Tucker, and Jasper Snoek. Deep Bayesian bandits showdown: An
 374 empirical comparison of Bayesian deep networks for Thompson sampling. *arXiv preprint
 375 arXiv:1802.09127*, 2018.
- 376 [44] Daniel Russo and Benjamin Van Roy. Learning to optimize via posterior sampling. *Mathematics
 377 of Operations Research*, 39(4):1221–1243, 2014.
- 378 [45] Daniel Russo, Benjamin Van Roy, Abbas Kazerouni, Ian Osband, and Zheng Wen. A tutorial
 379 on thompson sampling. *arXiv preprint arXiv:1707.02038*, 2017.
- 380 [46] Holger Schwenk and Yoshua Bengio. Boosting neural networks. *Neural computation*, 12(8):
 381 1869–1887, 2000.
- 382 [47] William R Thompson. On the likelihood that one unknown probability exceeds another in view
 383 of the evidence of two samples. *Biometrika*, 25(3/4):285–294, 1933.
- 384 [48] Michal Valko, Nathaniel Korda, Rémi Munos, Ilias Flaounas, and Nelo Cristianini. Finite-time
 385 analysis of kernelised contextual bandits. *arXiv preprint arXiv:1309.6869*, 2013.
- 386 [49] Lin F Yang and Mengdi Wang. Reinforcement leaning in feature space: Matrix bandit, kernels,
 387 and regret bound. *arXiv preprint arXiv:1905.10389*, 2019.
- 388 [50] Tom Zahavy and Shie Mannor. Deep neural linear bandits: Overcoming catastrophic forgetting
 389 through likelihood matching. *arXiv preprint arXiv:1901.08612*, 2019.
- 390 [51] Dongruo Zhou, Lihong Li, and Quanquan Gu. Neural contextual bandits with UCB-based
 391 exploration. *arXiv preprint arXiv:1911.04462*, 2019.
- 392 [52] Difan Zou, Yuan Cao, Dongruo Zhou, and Quanquan Gu. Gradient descent optimizes over-
 393 parameterized deep relu networks. *Machine Learning*, pages 1–26, 2019.

394 **A Further results of the experiments in Section 5**

395 Table 1 summarizes the total regrets measured at the last round on different data sets, with mean
 396 and standard deviation error computed based on 8 independent runs. The **Bold Faced** data is the
 397 top performance over 8 experiments. Table 2 shows the number of times the algorithm in that row
 398 significantly outperforms, ties, or significantly underperforms, compared with other algorithm with
 t -test at 90% significance level.

Table 1: Total regrets get at the last step with standard deviation attached

	Adult	Covertype	Magic ¹	MNIST	Mushroom	Shuttle
Round#	10 000	10 000	10 000	10 000	8 124	10 000
Input Dim ²	2×15	2×55	2×12	10×784	2×23	7×9
Random ³	5000	5000	5000	9000	4062	8571
Linear UCB	2078.0 ± 47.1	3220.4 ± 59.0	2616.2 ± 29.6	2544.0 ± 235.4	569.6 ± 18.1	956.5 ± 22.9
Linear TS	2118.1 ± 41.7	3385.4 ± 72.1	2605.2 ± 33.3	2781.4 ± 338.3	625.4 ± 60.7	1045.6 ± 53.8
Kernel UCB	2060.5 ± 20.1	3547.2 ± 103.9	2405.1 ± 85.6	3399.5 ± 258.4	182.9 ± 32.9	182.2 ± 24.3
Kernel TS	2110.2 ± 88.3	3693.0 ± 123.6	2415.9 ± 47.5	3385.1 ± 401.0	278.9 ± 37.6	270.2 ± 63.8
BootstrapNN	2095.5 ± 44.8	3060.2 ± 66.1	2267.2 ± 30.8	1776.6 ± 380.9	130.5 ± 9.9	210.6 ± 25.2
ϵ -greedy	2328.5 ± 50.4	3334.2 ± 72.6	2381.8 ± 37.3	1893.2 ± 93.7	323.2 ± 32.5	682.0 ± 79.8
Neural UCB	2102.8 ± 33.1	3058.5 ± 39.3	2074.0 ± 43.6	1531.0 ± 268.4	84.5 ± 23.7	209.6 ± 105.8
Neural TS (ours)	2088.2 ± 69.8	3069.2 ± 73.1	2088.4 ± 54.6	1522.8 ± 194.6	83.1 ± 37.4	242.5 ± 206.7

Table 2: Performance on total regret comparing with other methods on all datasets. Tuple ($w/t/l$) indicates the times of the algorithm at that row wins, ties with or loses, compared to all other 7 algorithms with t -test at 90% significant level.

	Adult	Covertype	Magic ⁴	MNIST	Mushroom	Shuttle
Linear UCB	1/6/0	4/0/3	0/1/6	2/1/4	1/0/6	1/0/6
Linear TS	1/5/1	2/1/4	0/1/6	2/1/4	0/0/7	0/0/7
Kernel UCB	4/3/0	1/0/6	2/2/3	0/1/6	4/0/3	5/2/0
Kernel TS	1/6/0	0/0/7	2/2/3	0/1/6	3/0/4	3/2/2
BootstrapNN	1/5/1	5/2/0	5/0/2	4/3/0	5/0/2	4/2/1
ϵ -greedy	0/0/7	2/1/4	2/2/3	4/1/2	2/0/5	2/0/5
Neural UCB	1/5/1	5/2/0	6/1/0	5/2/0	6/1/0	3/4/0
Neural TS (ours)	1/6/0	5/2/0	6/1/0	5/2/0	6/1/0	3/4/0

399

¹Magic is short for data set MagicTelescope

²Using disjoint encoding thus is NumofClass \times NumofFeatures

³Random pulling an arm at each round

⁴Magic is short for data set MagicTelescope

400 **B Proof of Lemmas in Section 4**

401 **B.1 Proof of Lemma 4.2**

402 The following concentration bound on Gaussian distributions will be useful in our proof.

403 **Lemma B.1** (Hoffman et al. [25]). Consider a normally distributed random variable $X \sim \mathcal{N}(\mu, \sigma^2)$ and $\beta \geq 0$. The probability that X is within a radius of $\beta\sigma$ from its mean can then be written as

$$\Pr(|X - \mu| \leq \beta\sigma) \geq 1 - \exp(-\beta^2/2).$$

405 *Proof of Lemma 4.2.* Since the estimated reward $\tilde{r}_{t,k}$ is sampled from $\mathcal{N}(f(\mathbf{x}_{t,k}; \boldsymbol{\theta}_{t-1}), \nu^2 \sigma_{t,k}^2)$ if
406 given filtration \mathcal{F}_t , Lemma B.1 implies that, conditioned on \mathcal{F}_t and given t, k ,

$$\Pr(|\tilde{r}_{t,k} - f(\mathbf{x}_{t,k}; \boldsymbol{\theta}_{t-1})| \leq c_t \nu \sigma_{t,k} | \mathcal{F}_t) \geq 1 - \exp(-c_t^2/2).$$

407 Taking a union bound over K arms, we have that for any t

$$\Pr(\forall k, |\tilde{r}_{t,k} - f(\mathbf{x}_{t,k}; \boldsymbol{\theta}_{t-1})| \leq c_t \nu \sigma_{t,k} | \mathcal{F}_t) \leq 1 - K \exp(-c_t^2/2).$$

408 Finally, choose $c_t = \sqrt{4 \log t + 2 \log K}$ as defined in (4.1), we get the bound that

$$\Pr(\mathcal{E}_t^\sigma | \mathcal{F}_t) = \Pr(\forall k, |\tilde{r}_{t,k} - f(\mathbf{x}_{t,k}; \boldsymbol{\theta}_{t-1})| \leq c_t \nu \sigma_{t,k} | \mathcal{F}_t) \leq 1 - \frac{1}{t^2}.$$

409 □

410 **B.2 Proof of Lemma 4.3**

411 Before going into the proof, some notation is needed about linear and kernelized models.

412 **Definition B.2.** Define $\bar{\mathbf{U}}_t = \lambda \mathbf{I} + \sum_{i=1}^t \mathbf{g}(\mathbf{x}_{i,a_i}; \boldsymbol{\theta}_0) \mathbf{g}(\mathbf{x}_{i,a_i}; \boldsymbol{\theta}_0)^\top / m$ and based on $\bar{\mathbf{U}}_t$, we further
413 define $\bar{\sigma}_{t,k}^2 = \lambda \mathbf{g}^\top(\mathbf{x}_{t,k}; \boldsymbol{\theta}_0) \bar{\mathbf{U}}_{t-1}^{-1} \mathbf{g}(\mathbf{x}_{t,k}; \boldsymbol{\theta}_0) / m$. Furthermore, for convenience we define

$$\begin{aligned} \mathbf{J}_t &= (\mathbf{g}(\mathbf{x}_{1,a_1}; \boldsymbol{\theta}_t) \quad \cdots \quad \mathbf{g}(\mathbf{x}_{t,a_t}; \boldsymbol{\theta}_t)), \\ \bar{\mathbf{J}}_t &= (\mathbf{g}(\mathbf{x}_{1,a_1}; \boldsymbol{\theta}_0) \quad \cdots \quad \mathbf{g}(\mathbf{x}_{t,a_t}; \boldsymbol{\theta}_0)), \\ \mathbf{h}_t &= (h(\mathbf{x}_{1,a_1}) \quad \cdots \quad h(\mathbf{x}_{t,a_t}))^\top, \\ \mathbf{r}_t &= (r_1 \quad \cdots \quad r_t)^\top, \\ \boldsymbol{\epsilon}_t &= (h(\mathbf{x}_{1,a_1}) - r_1 \quad \cdots \quad h(\mathbf{x}_{t,a_t}) - r_t)^\top, \end{aligned}$$

414 where $\boldsymbol{\epsilon}_t$ is the reward noise. We can verify that $\mathbf{U}_t = \lambda \mathbf{I} + \mathbf{J}_t \mathbf{J}_t^\top / m$, $\bar{\mathbf{U}}_t = \lambda \mathbf{I} + \bar{\mathbf{J}}_t \bar{\mathbf{J}}_t^\top / m$, and
415 $\mathbf{h}_t = \bar{\mathbf{J}}_t^\top(\boldsymbol{\theta}^* - \boldsymbol{\theta}_0)$. We further define $\mathbf{K}_t = \bar{\mathbf{J}}_t^\top \bar{\mathbf{J}}_t / m$.

416 The first lemma shows that the target function is well-approximated by the linearized neural network
417 if the network width m is large enough.

418 **Lemma B.3** (Lemma 5.1, Zhou et al. [51]). There exists some constant $C > 0$ such that for any
419 $\delta \in (0, 1)$, if

$$m \geq CT^4 K^4 L^6 \log(T^2 K^2 L / \delta) / \lambda_0^4,$$

420 then with probability at least $1 - \delta$ over the random initialization of $\boldsymbol{\theta}_0$, there exists a $\boldsymbol{\theta}^* \in \mathbb{R}^p$ such
421 that

$$h(\mathbf{x}^i) = \langle \mathbf{g}(\mathbf{x}^i; \boldsymbol{\theta}_0), \boldsymbol{\theta}^* - \boldsymbol{\theta}_0 \rangle, \quad \sqrt{m} \|\boldsymbol{\theta}^* - \boldsymbol{\theta}_0\|_2 \leq \sqrt{2 \mathbf{h}^\top \mathbf{H}^{-1} \mathbf{h}} \leq B, \quad (\text{B.1})$$

422 for all $i \in [TK]$, where B is defined in Theorem 3.5.

423 The next lemma bounds the difference between the $\bar{\sigma}_{t,k}$ from the linearized model and the $\sigma_{t,k}$
424 actually used in the algorithm. Its proof, together with other technical lemmas', will be given in the
425 next section.

426 **Lemma B.4.** If the network size m and the learning rate η satisfy the condition in Condition 4.1,
427 then with probability at least $1 - \delta$,

$$|\bar{\sigma}_{t,k} - \sigma_{t,k}| \leq C\sqrt{\log m} t^{7/6} m^{-1/6} \lambda^{-2/3} L^{9/2},$$

428 where C is a positive constant.

429 We next bound the difference between the outputs of the neural network and the linearized model.

430 **Lemma B.5.** Suppose the network width m and learning rate η satisfy the condition in Condition 4.1.
431 Then, with probability at least $1 - \delta$ over the random initialization of θ_0 , we have

$$\begin{aligned} |f(\mathbf{x}_{t,k}) - \langle \mathbf{g}(x_{t,k}; \theta_0), \bar{\mathbf{U}}_{t-1}^{-1} \bar{\mathbf{J}}_{t-1} \mathbf{r}_{t-1} / m \rangle| &\leq C_1 t^{2/3} m^{-1/6} \lambda^{-2/3} L^3 \sqrt{\log m} \\ &\quad + C_2 (1 - \eta m \lambda)^J \sqrt{tL/\lambda} \\ &\quad + \bar{C}_3 m^{-1/6} \sqrt{\log m} L^4 t^{5/3} \lambda^{-5/3} (1 + \sqrt{t/\lambda}), \end{aligned}$$

432 where $\{C_i\}_{i=1}^3$ are positive constants.

433 The next lemma, due to Chowdhury and Gopalan [15], controls the quadratic value generated by an
434 R -sub-Gaussian random vector ϵ :

435 **Lemma B.6** (Theorem 1, Chowdhury and Gopalan [15]). Let $\{\epsilon_t\}_{t=1}^\infty$ be a real-valued stochastic
436 process such that for some $R \geq 0$ and for all $t \geq 1$, ϵ_t is \mathcal{F}_t -measurable and R -sub-Gaussian
437 conditioned on \mathcal{F}_t . Recall \mathbf{K}_t defined in Definition B.2. With probability $0 < \delta < 1$ and for a given
438 $\eta > 0$, with probability $1 - \delta$, the following holds for all t ,

$$\epsilon_{1:t}^\top ((\mathbf{K}_t + \eta \mathbf{I})^{-1} + \mathbf{I})^{-1} \epsilon_{1:t} \leq R^2 \log \det((1 + \eta) \mathbf{I} + \mathbf{K}_t) + 2R^2 \log(1/\delta).$$

439 Finally, the following lemma shows the linearized kernel and the neural tangent kernel are closed:

440 **Lemma B.7.** For all $t \in [T]$, there exists a positive constants C such that the following holds: if the
441 network width m satisfies

$$m \geq CT^6 L^6 K^6 \log(TKL/\delta),$$

442 then with probability at least $1 - \delta$,

$$\log \det(\mathbf{I} + \lambda^{-1} \mathbf{K}_t) \leq \log \det(\mathbf{I} + \lambda^{-1} \mathbf{H}) + 1.$$

443 We are now ready to prove Lemma 4.3.

444 *Proof of Lemma 4.3.* First of all, it is easy to verify that the condition of η and m on Condition 4.1
445 satisfies the condition required in Lemmas B.3–B.7. Thus, taking a union bound, we
446 have with probability at least $1 - 5\delta$, that the bounds provided by these lemmas hold. Then for
447 any $t \in [T]$, we will first provide the difference between the target function and the linear function
448 $\langle \mathbf{g}(x_{t,k}; \theta_0), \bar{\mathbf{U}}_{t-1}^{-1} \bar{\mathbf{J}}_{t-1} \mathbf{r}_{t-1} / m \rangle$ as:

$$\begin{aligned} &|h(\mathbf{x}_{t,k}) - \langle \mathbf{g}(x_{t,k}; \theta_0), \bar{\mathbf{U}}_{t-1}^{-1} \bar{\mathbf{J}}_{t-1} \mathbf{r}_{t-1} / m \rangle| \\ &\leq |h(\mathbf{x}_{t,k}) - \langle \mathbf{g}(x_{t,k}; \theta_0), \bar{\mathbf{U}}_{t-1}^{-1} \bar{\mathbf{J}}_{t-1} \mathbf{h}_{t-1} / m \rangle| + |\langle \mathbf{g}(x_{t,k}; \theta_0), \bar{\mathbf{U}}_{t-1}^{-1} \bar{\mathbf{J}}_{t-1} \epsilon_{t-1} / m \rangle| \quad (\text{B.2}) \\ &= |\langle \mathbf{g}(x_{t,k}; \theta_0), \theta^* - \theta_0 - \bar{\mathbf{U}}_{t-1}^{-1} \bar{\mathbf{J}}_{t-1}^\top (\theta^* - \theta_0) / m \rangle| + |\mathbf{g}(x_{t,k})^\top \bar{\mathbf{U}}_{t-1}^{-1} \bar{\mathbf{J}}_{t-1} \epsilon_{t-1} / m| \quad (\text{B.3}) \end{aligned}$$

$$= |\langle \mathbf{g}(x_{t,k}; \theta_0), (\mathbf{I} - \bar{\mathbf{U}}_{t-1}^{-1} (\bar{\mathbf{U}}_{t-1} - \lambda \mathbf{I})) (\theta^* - \theta_0) \rangle| + |\mathbf{g}(x_{t,k}; \theta_0)^\top \bar{\mathbf{U}}_{t-1}^{-1} \bar{\mathbf{J}}_{t-1} \epsilon_{t-1} / m| \quad (\text{B.4})$$

$$\begin{aligned} &= \lambda |\mathbf{g}(x_{t,k}; \theta_0)^\top \bar{\mathbf{U}}_{t-1}^{-1} (\theta^* - \theta_0)| + |\mathbf{g}(x_{t,k}; \theta_0)^\top \bar{\mathbf{U}}_{t-1}^{-1} \bar{\mathbf{J}}_{t-1} \epsilon_{t-1} / m| \\ &\leq \lambda \sqrt{\mathbf{g}(x_{t,k}; \theta_0)^\top \bar{\mathbf{U}}_{t-1}^{-1} \mathbf{g}(x_{t,k}; \theta_0)} \sqrt{(\theta^* - \theta_0)^\top \bar{\mathbf{U}}_{t-1}^{-1} (\theta^* - \theta_0)} \\ &\quad + \sqrt{\mathbf{g}(x_{t,k}; \theta_0)^\top \bar{\mathbf{U}}_{t-1}^{-1} \mathbf{g}(x_{t,k}; \theta_0)} \sqrt{\epsilon_{t-1}^\top \bar{\mathbf{J}}_{t-1}^\top \bar{\mathbf{U}}_{t-1}^{-1} \bar{\mathbf{J}}_{t-1} \epsilon_{t-1} / m} \quad (\text{B.5}) \end{aligned}$$

$$\leq \sqrt{m} \|\theta^* - \theta_0\|_2 \bar{\sigma}_{t,k} + \bar{\sigma}_{t,k} \lambda^{-1/2} \sqrt{\epsilon_{t-1}^\top \bar{\mathbf{J}}_{t-1}^\top \bar{\mathbf{U}}_{t-1}^{-1} \bar{\mathbf{J}}_{t-1} \epsilon_{t-1} / m} \quad (\text{B.6})$$

449 where: (B.2) uses triangle inequality and the fact that $\mathbf{r}_{t-1} = \mathbf{h}_{t-1} + \epsilon_{t-1}$; (B.3) is from Lemma B.3;
450 (B.4) uses the fact that $\bar{\mathbf{J}}_{t-1}^\top \bar{\mathbf{J}}_{t-1} = m(\bar{\mathbf{U}}_{t-1} - \lambda \mathbf{I})$ which can be verified using Definition B.2;

451 (B.5) is from the fact that $|\boldsymbol{\alpha}^\top \mathbf{A} \boldsymbol{\beta}| \leq \sqrt{\boldsymbol{\alpha}^\top \mathbf{A} \boldsymbol{\alpha}} \sqrt{\boldsymbol{\beta}^\top \mathbf{A} \boldsymbol{\beta}}$. Since $\mathbf{U}_{t-1}^{-1} \preceq \frac{1}{\lambda} \mathbf{I}$ and $\bar{\sigma}_{t,k}$ defined in
452 Definition B.2, we obtain (B.6).

453 Furthermore, by obtaining

$$\begin{aligned}\bar{\mathbf{J}}_{t-1}^\top \mathbf{U}_{t-1}^{-1} \bar{\mathbf{J}}_{t-1} / m &= \bar{\mathbf{J}}_{t-1}^\top (\lambda \mathbf{I} + \bar{\mathbf{J}}_{t-1} \bar{\mathbf{J}}_{t-1}^\top / m)^{-1} \bar{\mathbf{J}}_{t-1} \\ &= \bar{\mathbf{J}}_{t-1}^\top (\lambda^{-1} \mathbf{I} - \lambda^{-2} \bar{\mathbf{J}}_{t-1} (\mathbf{I} + \lambda^{-1} \bar{\mathbf{J}}_{t-1}^\top \bar{\mathbf{J}}_{t-1} / m)^{-1} \bar{\mathbf{J}}_{t-1}^\top / m) \bar{\mathbf{J}}_{t-1} / m \quad (\text{B.7}) \\ &= \lambda^{-1} \bar{\mathbf{J}}_{t-1}^\top \bar{\mathbf{J}}_{t-1} / m - \lambda^{-1} \bar{\mathbf{J}}_{t-1}^\top \bar{\mathbf{J}}_{t-1} (\lambda \mathbf{I} + \bar{\mathbf{J}}_{t-1}^\top \bar{\mathbf{J}}_{t-1} / m)^{-1} \bar{\mathbf{J}}_{t-1}^\top \bar{\mathbf{J}}_{t-1} / m^2 \\ &= \mathbf{K}_{t-1} (\lambda^{-1} \mathbf{I} + (\lambda \mathbf{I} + \mathbf{K}_{t-1})^{-1} \mathbf{K}_{t-1}) = \mathbf{K}_{t-1} (\lambda \mathbf{I} + \mathbf{K}_{t-1})^{-1}, \quad (\text{B.8})\end{aligned}$$

454 where (B.7) is from the Sherman-Morrison formula, and (B.8) uses Definition B.2 and the fact that
455 $(\lambda^{-1} \mathbf{I} + \mathbf{K}_{t-1})^{-1} \mathbf{K}_{t-1} = \lambda \mathbf{I} + \mathbf{K}_{t-1}$ which could be verified by multiplying the LHS and RHS
456 together, we have that

$$\sqrt{\epsilon_{t-1}^\top \bar{\mathbf{J}}_{t-1}^\top \mathbf{U}_{t-1}^{-1} \bar{\mathbf{J}}_{t-1} \epsilon_{t-1}} / m \leq \sqrt{\epsilon_{t-1}^\top \mathbf{K}_{t-1} (\lambda \mathbf{I} + \mathbf{K}_{t-1})^{-1} \epsilon_{t-1}} \quad (\text{B.9})$$

$$\leq \sqrt{\epsilon_{t-1}^\top (\mathbf{K}_{t-1} + (\lambda - 1) \mathbf{I}) (\lambda \mathbf{I} + \mathbf{K}_{t-1})^{-1} \epsilon_{t-1}} \quad (\text{B.10})$$

$$= \sqrt{\epsilon_{t-1}^\top (\mathbf{I} + (\mathbf{K}_{t-1} + (\lambda - 1) \mathbf{I})^{-1})^{-1} \epsilon_{t-1}} \quad (\text{B.11})$$

457 where (B.10) is because $\lambda = 1 + 1/T \geq 1$ set in Theorem 3.5.

458 Based on (B.6) and (B.11), by utilizing the bound on $\|\boldsymbol{\theta}^* - \boldsymbol{\theta}\|_2$ provided in Lemma B.3, as well as
459 the bound given in Lemma B.6, and $\lambda \geq 1$, we have

$$|h(\mathbf{x}_{t,k}) - \langle \mathbf{g}(\mathbf{x}_{t,k}; \boldsymbol{\theta}_0), \bar{\mathbf{U}}_{t-1}^{-1} \bar{\mathbf{J}}_{t-1} \mathbf{r}_{t-1} m \rangle| \leq (B + R \sqrt{\log \det(\lambda \mathbf{I} + \mathbf{K}_{t-1}) + 2 \log(1/\delta)}) \bar{\sigma}_{t,k},$$

460 since it is obvious that

$$\begin{aligned}\log \det(\lambda \mathbf{I} + \mathbf{K}_{t-1}) &= \log \det(\mathbf{I} + \lambda^{-1} \mathbf{K}_{t-1}) + (t-1) \log \lambda \\ &\leq \log \det(\mathbf{I} + \lambda^{-1} \mathbf{K}_{t-1}) + t(\lambda - 1) \\ &\leq \log \det(\mathbf{I} + \lambda^{-1} \mathbf{H}) + 2,\end{aligned}$$

461 where the first equality moves the λ outside the $\log \det$, the first inequality is due to $\log \lambda \leq \lambda - 1$,
462 and the second inequality is from Lemma B.7 and the fact that $\lambda = 1 + 1/T$ (as set in Theorem 3.5).
463 Thus, we have

$$|h(\mathbf{x}_{t,k}) - \langle \mathbf{g}(\mathbf{x}_{t,k}; \boldsymbol{\theta}_0), \bar{\mathbf{U}}_{t-1}^{-1} \bar{\mathbf{J}}_{t-1} \mathbf{r}_{t-1} m \rangle| \leq \nu \bar{\sigma}_{t,k},$$

464 where we set $\nu = B + R \sqrt{\log \det(\mathbf{I} + \mathbf{H}/\lambda) + 2 + 2 \log(1/\delta)}$. Then, by combining this bound
465 with Lemma B.5, we conclude that there exist positive constants $\bar{C}_1, \bar{C}_2, \bar{C}_3$ so that

$$\begin{aligned}|f(\mathbf{x}_{t,k}) - h(\mathbf{x}_{t,k})| &\leq \nu \bar{\sigma}_{t,k} + \bar{C}_1 t^{2/3} m^{-1/6} \lambda^{-2/3} L^3 \sqrt{\log m} + \bar{C}_2 (1 - \eta m \lambda)^J \sqrt{tL/\lambda} \\ &\quad + \bar{C}_3 m^{-1/6} \sqrt{\log m} L^4 t^{5/3} \lambda^{-5/3} (1 + \sqrt{t/\lambda}), \\ &\leq \nu \bar{\sigma}_{t,k} + \bar{C}_1 t^{2/3} m^{-1/6} \lambda^{-2/3} L^3 \sqrt{\log m} + \bar{C}_2 (1 - \eta m \lambda)^J \sqrt{tL/\lambda} \\ &\quad + \bar{C}_3 m^{-1/6} \sqrt{\log m} L^4 t^{5/3} \lambda^{-5/3} (1 + \sqrt{t/\lambda}) \\ &\quad + (B + R \sqrt{\log \det(\mathbf{I} + \mathbf{H}/\lambda) + 2 + 2 \log(1/\delta)}) (\bar{\sigma}_{t,k} - \sigma_{t,k}).\end{aligned}$$

466 Finally, by utilizing the bound of $|\bar{\sigma}_{t,k} - \sigma_{t,k}|$ provided in Lemma B.4, we conclude that

$$|f(\mathbf{x}_{t,k}) - h(\mathbf{x}_{t,k})| \leq \nu \sigma_{t,k} + \delta_m,$$

467 where $\epsilon(m)$ is defined by adding all of the additional terms and taking $t = T$:

$$\begin{aligned}\epsilon(m) &= \bar{C}_1 T^{2/3} m^{-1/6} \lambda^{-2/3} L^3 \sqrt{\log m} + \bar{C}_2 (1 - \eta m \lambda)^J \sqrt{TL/\lambda} + \\ &\quad + \bar{C}_3 m^{-1/6} \sqrt{\log m} L^4 t^{5/3} \lambda^{-5/3} (1 + \sqrt{T/\lambda}) \\ &\quad + \bar{C}_4 (B + R \sqrt{\log \det(\mathbf{I} + \mathbf{H}/\lambda) + 2 + 2 \log(1/\delta)}) \sqrt{\log m} T^{7/6} m^{-1/6} \lambda^{-2/3} L^{9/2},\end{aligned}$$

468 where is exactly the same form defined in (4.3). By setting δ to $\delta/5$ (required by the union bound
469 discussed at the beginning of the proof), we get the result presented in Lemma 4.3. \square

470 **B.3 Proof of Lemma 4.4**

471 Our proof requires an anti-concentration bound for Gaussian distribution, as stated below:

472 **Lemma B.8** (Gaussian anti-concentration). For a Gaussian random variable X with mean μ and
473 standard deviation σ , for any $\beta > 0$,

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

474 *Proof of Lemma 4.4.* Since $\tilde{r}_{t,k} \sim \mathcal{N}(f(\mathbf{x}_{t,k}; \boldsymbol{\theta}_{t-1}), \nu_t^2 \sigma_{t,k}^2)$ conditioned on \mathcal{F}_t , we have

$$\begin{aligned} & \Pr(\tilde{r}_{t,k} + \epsilon(m) > h(\mathbf{x}_{t,k}) \mid \mathcal{F}_t, \mathcal{E}_t^\mu) \\ &= \Pr\left(\frac{\tilde{r}_{t,k} - f(\mathbf{x}_{t,k}; \boldsymbol{\theta}_{t-1}) + \epsilon(m)}{\nu \sigma_{t,k}} > \frac{h(x_{t,k}) - f(\mathbf{x}_{t,k}; \boldsymbol{\theta}_{t-1})}{\nu \sigma_{t,k}} \mid \mathcal{F}_t, \mathcal{E}_t^\mu\right) \\ &\geq \Pr\left(\frac{\tilde{r}_{t,k} - f(\mathbf{x}_{t,k}; \boldsymbol{\theta}_{t-1}) + \epsilon(m)}{\nu \sigma_{t,k}} > \frac{|h(x_{t,k}) - f(\mathbf{x}_{t,k}; \boldsymbol{\theta}_{t-1})|}{\nu \sigma_{t,k}} \mid \mathcal{F}_t, \mathcal{E}_t^\mu\right) \\ &= \Pr\left(\frac{\tilde{r}_{t,k} - f(\mathbf{x}_{t,k}; \boldsymbol{\theta}_{t-1})}{\nu \sigma_{t,k}} > \frac{|h(x_{t,k}) - f(\mathbf{x}_{t,k}; \boldsymbol{\theta}_{t-1})| - \epsilon(m)}{\nu \sigma_{t,k}} \mid \mathcal{F}_t, \mathcal{E}_t^\mu\right) \\ &\geq \Pr\left(\frac{\tilde{r}_{t,k} - f(\mathbf{x}_{t,k}; \boldsymbol{\theta}_{t-1})}{\nu \sigma_{t,k}} > 1 \mid \mathcal{F}_t, \mathcal{E}_t^\mu\right) \geq \frac{1}{4e\sqrt{\pi}}, \end{aligned}$$

475 where the first inequality is due to $|x| \geq x$, and the second inequality follows from event \mathcal{E}_t^μ , i.e.,

$$\forall k \in [K], \quad |f(\mathbf{x}_{t,k}; \boldsymbol{\theta}_{t-1}) - h(\mathbf{x}_{t,k})| \leq \nu \sigma_{t,k} + \epsilon(m).$$

476 \square

477 **B.4 Proof of Lemma 4.5**

478 *Proof of Lemma 4.5.* Consider the following two events at round t :

$$\begin{aligned} \mathcal{A} &= \{\forall k \in S_t, \tilde{r}_{t,k} < \tilde{r}_{t,a_t^*} \mid \mathcal{F}_t, \mathcal{E}_t^\mu\}, \\ \mathcal{B} &= \{a_t \notin S_t \mid \mathcal{F}_t, \mathcal{E}_t^\mu\}. \end{aligned}$$

479 Clearly, \mathcal{A} implies \mathcal{B} , since $a_t = \operatorname{argmax}_k \tilde{r}_{t,k}$. Therefore,

$$\Pr(a_t \notin S_t \mid \mathcal{F}_t, \mathcal{E}_t^\mu) \geq \Pr(\forall k \in S_t, \tilde{r}_{t,k} < \tilde{r}_{t,a_t^*} \mid \mathcal{F}_t, \mathcal{E}_t^\mu).$$

480 Suppose \mathcal{E}^θ also holds, then it is easy to show that $\forall k \in [K]$,

$$|h(\mathbf{x}_{t,k}) - \tilde{r}_{t,k}| \leq |h(\mathbf{x}_{t,k}) - f(\mathbf{x}_{t,k}; \boldsymbol{\theta}_t)| + |f(\mathbf{x}_{t,k}; \boldsymbol{\theta}_t) - \tilde{r}_{t,k}| \leq \epsilon(m) + (1 + \tilde{c}_t)\nu_t \sigma_{t,k}. \quad (\text{B.12})$$

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

$$h(\mathbf{x}_{t,a_t^*}) - \tilde{r}_{t,k} \geq h(\mathbf{x}_{t,a_t^*}) - h(\mathbf{x}_{t,k}) - |h(\mathbf{x}_{t,k}) - \tilde{r}_{t,k}| \geq \epsilon(m),$$

482 where we used the definitions of saturated arms in Definition 4.4, and of \mathcal{E}_t^μ and \mathcal{E}_t^σ in 4.1.

Consider the following event

$$\mathcal{C} = \{h(\mathbf{x}_{t,a_t^*}) - \epsilon(m) < \tilde{r}_{t,a_t^*} \mid \mathcal{F}_t, \mathcal{E}_t^\mu\}.$$

483 Since \mathcal{E}_t^σ implies $h(\mathbf{x}_{t,a_t^*}) - \epsilon(m) \geq \tilde{r}_{t,k}$, we have that if $\mathcal{C}, \mathcal{E}_t^\sigma$ holds, then \mathcal{A} holds, i.e. $\mathcal{E}_t^\sigma \cap \mathcal{C} \subseteq \mathcal{A}$.

484 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

$$\Pr(\mathcal{A}) + \Pr(\bar{\mathcal{E}}_t^\sigma) \geq \Pr(\mathcal{C}). \quad (\text{B.13})$$

485 Then, (B.13) implies that

$$\begin{aligned} \Pr(\forall k \in S_t, \tilde{r}_{t,k} < \tilde{r}_{t,a_t^*} \mid \mathcal{F}_t, \mathcal{E}_t^\mu) &\geq \Pr(\tilde{r}_{t,a_t^*} + \epsilon(m) > h(\mathbf{x}_{t,a_t^*}) \mid \mathcal{F}_t, \mathcal{E}_t^\mu) - \Pr(\bar{\mathcal{E}}_t^\sigma \mid \mathcal{F}_t, \mathcal{E}_t^\mu) \\ &\geq \frac{1}{4e\sqrt{\pi}} - \frac{1}{t^2}, \end{aligned}$$

486 where the first inequality is from a_t^* is a special case of $\forall k \in [K]$, the second inequality is from
487 Lemmas 4.2 and 4.4. \square

488 **B.5 Proof of Lemma 4.6**

489 To prove Lemma 4.6, we will need an upper bound bound on $\delta_{t,k}$.

490 **Lemma B.9.** For any time $t \in [T]$, $k \in [K]$, and $\delta \in (0, 1)$, if the learning rate η and network width
491 m satisfy Condition 4.1, we have, with probability at least $1 - \delta$, that

$$\sigma_{t,k} \leq C\sqrt{L},$$

492 where C is a positive constant.

493 *Proof of Lemma 4.6.* Recall that given \mathcal{F}_t and \mathcal{E}_t^μ , the only randomness comes from sampling $\tilde{r}_{t,k}$
494 for $k \in [K]$. Let \bar{k}_t be the unsaturated arm with the smallest $\sigma_{t,..}$, i.e.

$$\bar{k}_t = \operatorname{argmin}_{k \notin S_t} \sigma_{t,k},$$

495 then we have that

$$\mathbb{E}[\sigma_{t,a_t} | \mathcal{F}_t, \mathcal{E}_t^\mu] \geq \mathbb{E}[\sigma_{t,a_t} | \mathcal{F}_t, \mathcal{E}_t^\mu, a_t \notin S_t] \Pr(a_t \notin S_t | \mathcal{F}_t, \mathcal{E}_t^\mu) \quad (\text{B.14})$$

$$\geq \sigma_{t,\bar{k}_t} \left(\frac{1}{4e\sqrt{\pi}} - \frac{1}{t^2} \right), \quad (\text{B.15})$$

496 where the first inequality ignores the case when $a_t \in S_t$, and the second inequality is from Lemma 4.5
497 and the definition of \bar{k}_t mentioned above. If both \mathcal{E}_t^σ and \mathcal{E}_t^μ hold, then

$$\forall k \in [K], |h(\mathbf{x}_{t,k}) - \tilde{r}_{t,k}| \leq \epsilon(m) + (1 + c_t)\nu\sigma_{t,k}, \quad (\text{B.16})$$

498 as proved in equation (B.12). Thus,

$$\begin{aligned} h(\mathbf{x}_{t,a_t^*}) - h(\mathbf{x}_{t,a_t}) &= h(\mathbf{x}_{t,a_t^*}) - h(\mathbf{x}_{t,\bar{k}_t}) + h(\mathbf{x}_{t,\bar{k}_t}) - h(\mathbf{x}_{t,a_t}) \\ &\leq (1 + c_t)\nu\sigma_{t,\bar{k}_t}, 1 + 2\epsilon(m) + h(\mathbf{x}_{t,\bar{k}_t}) - \tilde{r}_{t,\bar{k}_t} - h(\mathbf{x}_{t,a_t}) \end{aligned} \quad (\text{B.17})$$

$$\begin{aligned} &\quad + \tilde{r}_{t,a_t} + \tilde{r}_{t,\bar{k}_t} - \tilde{r}_{t,a_t} \\ &\leq (1 + c_t)\nu(2\sigma_{t,\bar{k}_t} + \sigma_{t,a_t}) + 4\epsilon(m), \end{aligned} \quad (\text{B.18})$$

499 where the first inequality is from Definition 4.4 and $\bar{k}_t \notin S_t$, and the second inequality comes from
500 equation (B.16). Since a trivial bound on $h(\mathbf{x}_{t,a_t^*}) - h(\mathbf{x}_{t,a_t})$ could be get by $h(\mathbf{x}_{t,a_t^*}) - h(\mathbf{x}_{t,a_t}) \leq$
501 $|h(\mathbf{x}_{t,a_t^*})| + |h(\mathbf{x}_{t,a_t})| \leq 2$, then we have

$$\begin{aligned} \mathbb{E}[h(\mathbf{x}_{t,a_t^*}) - h(\mathbf{x}_{t,a_t}) | \mathcal{F}_t, \mathcal{E}_t^\mu] &= \mathbb{E}[h(\mathbf{x}_{t,a_t^*}) - h(\mathbf{x}_{t,a_t}) | \mathcal{F}_t, \mathcal{E}_t^\mu, \mathcal{E}_t^\sigma] \Pr(\mathcal{E}_t^\sigma) \\ &\quad + \mathbb{E}[h(\mathbf{x}_{t,a_t^*}) - h(\mathbf{x}_{t,a_t}) | \mathcal{F}_t, \mathcal{E}_t^\mu, \bar{\mathcal{E}}_t^\sigma] \Pr(\bar{\mathcal{E}}_t^\sigma) \\ &\leq (1 + c_t)\nu(2\sigma_{t,\bar{k}_t} + \mathbb{E}[\sigma_{t,a_t} | \mathcal{F}_t, \mathcal{E}_t^\mu]) + 4\epsilon(m) + \frac{2}{t^2} \\ &\leq (1 + c_t)\nu \left(\frac{2\mathbb{E}[\sigma_{t,a_t} | \mathcal{F}_t, \mathcal{E}_t^\mu]}{\frac{1}{4e\sqrt{\pi}} - \frac{1}{t^2}} + \mathbb{E}[\sigma_{t,a_t} | \mathcal{F}_t, \mathcal{E}_t^\mu] \right) + 4\epsilon(m) + \frac{2}{t^2} \\ &\leq 44e\sqrt{\pi}(1 + c_t)\nu \mathbb{E}[\sigma_{t,a_t} | \mathcal{F}_t, \mathcal{E}_t^\mu] + 4\epsilon(m) + 2t^{-2}, \end{aligned}$$

502 where the inequality on the second line uses the bound provide in (B.18) and the trivial bound of
503 $h(\mathbf{x}_{t,a_t^*}) - h(\mathbf{x}_{t,a_t})$ for the second term plus Lemma 4.2, the inequality on the third line uses the
504 bound of σ_{t,\bar{k}_t} provide in (B.15), inequality on the forth line is directly calculated by $1 \leq 4e\sqrt{\pi}$ and

$$\frac{1}{\frac{1}{4e\sqrt{\pi}} - \frac{1}{t^2}} \leq 20e\sqrt{\pi},$$

505 which trivially holds since LHS is negative when $t \leq 4$ and when $t = 5$, the LHS reach its maximum
506 as $\approx 84.11 < 96.36 \approx$ RHS. Noticing that $|h(\mathbf{x})| \leq 1$, it is trivial to further extend the bound as

$$\mathbb{E}[h(\mathbf{x}_{t,a_t^*}) - h(\mathbf{x}_{t,a_t}) | \mathcal{F}_t, \mathcal{E}_t^\mu] \leq \min\{44e\sqrt{\pi}(1 + c_t)\nu \mathbb{E}[\sigma_{t,a_t} | \mathcal{F}_t, \mathcal{E}_t^\mu], 2\} + 4\epsilon(m) + 2t^{-2},$$

507 and since we have $1 + c_t \geq 1$ and $\nu = B + R\sqrt{\log \det(\mathbf{I} + \mathbf{H}/\lambda) + 2 + 2\log(1/\delta)} \geq B$, recall
508 $22e\sqrt{\pi}B \geq 1$, it is easy to verify the following inequality also holds:

$$\begin{aligned} & \mathbb{E}[h(\mathbf{x}_{t,a_t^*}) - h(\mathbf{x}_{t,a_t}) | \mathcal{F}_t, \mathcal{E}_t^\mu] \\ & \leq 44e\sqrt{\pi}(1 + c_t)\nu \min\{\mathbb{E}[\sigma_{t,a_t} | \mathcal{F}_t, \mathcal{E}_t^\mu], 1\} + 4\epsilon(m) + 2t^{-2} \\ & \leq 44e\sqrt{\pi}(1 + c_t)\nu C_1 \sqrt{L} \mathbb{E}[\min\{\sigma_{t,a_t}, 1\} | \mathcal{F}_t, \mathcal{E}_t^\mu] + 4\epsilon(m) + 2t^{-2}, \end{aligned}$$

509 where we use the fact that there exists a constant C_1 such that σ_{t,a_t} is bounded by $C_1\sqrt{L}$ with
510 probability $1 - \delta$ provided by Lemma B.9. Merging the positive constant C_1 with $44e\sqrt{\pi}$, we get
511 the statement in Lemma 4.6. \square

512 B.6 Proof of Lemma 4.7

513 We start with introducing the Azuma-Hoeffding inequality for super-martingale:

514 **Lemma B.10** (Azuma-Hoeffding Inequality for Super Martingale). If a super-martingale Y_t , cor-
515 responding to filtration \mathcal{F}_t satisfies that $|Y_t - Y_{t-1}| \leq B_t$, then for any $\delta \in (0, 1)$, w.p. $1 - \delta$, we
516 have

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

517 *Proof of Lemma 4.7.* With the condition of m and η described in Condition 4.1. From Lemma B.9,
518 we have that there exists a positive constant C_1 such that X_t defined in (4.5) is bounded with
519 probability $1 - \delta$ by

$$\begin{aligned} |X_t| & \leq |\bar{\Delta}_t| + C_1(1 + c_t)\nu\sqrt{L} \min\{\sigma_{t,a_t}, 1\} + 4\epsilon(m) + 2t^{-2} \\ & \leq 2 + 2t^{-2} + C_1C_2(1 + c_t)\nu L + 4\epsilon(m) \\ & \leq 4 + C_1C_2(1 + c_t)\nu L + 4\epsilon(m) \end{aligned}$$

520 where the first inequality uses the fact that $|a - b| \leq |a| + |b|$; the second inequality is from Lemma B.9
521 and the fact that $h \leq 1$, where C_2 is a positive constant used in Lemma B.9; the third inequality
522 uses the fact that $t^{-2} \leq 1$. Noticing the fact that $c_t \leq c_T$, and from Lemma 4.6 we know that with
523 probability at least $1 - \delta$, Y_t is a super martingale. From Lemma B.10, we have

$$Y_T - Y_0 \leq (4 + C_1C_2(1 + c_T)\nu L + 4\epsilon(m))\sqrt{2 \log(1/\delta)T}. \quad (\text{B.19})$$

524 Considering the definition of Y_T in (4.5), (B.19) is equivalent to

$$\begin{aligned} \sum_{i=1}^T \bar{\Delta}_i & \leq 4T\epsilon(m) + 2 \sum_{i=1}^T t^{-2} + C_1(1 + c_T)\nu\sqrt{L} \sum_{i=1}^T \min\{\sigma_{t,a_t}, 1\} \\ & \quad + (4 + C_1C_2(1 + c_T)\nu L + 4\epsilon(m))\sqrt{2 \log(1/\delta)T}, \end{aligned}$$

525 then by utilizing $\sum_{i=1}^\infty t^{-2} = \pi^2/6$, and merge the constant C_1 with $44e\sqrt{\pi}$, taking union bound of
526 the probability bound of Lemma 4.6, B.10, B.9 we have the inequality above hold with probability at
527 least $1 - 3\delta$. Re-scaling δ to $\delta/3$ and merging the product of C_1C_2 as a new positive constant leads
528 to the desired result. \square

529 B.7 Proof of Lemma 4.8

530 We first state a technical lemma that will be useful:

531 **Lemma B.11** (Lemma 11, Abbasi-Yadkori et al. [1]). Let $\{\mathbf{v}_t\}_{t=1}^\infty$ be a sequence in \mathbb{R}^d , and define
532 $\mathbf{V}_t = \lambda\mathbf{I} + \sum_{i=1}^t \mathbf{v}_i \mathbf{v}_i^\top$. If $\lambda \geq 1$, then

$$\sum_{i=1}^T \min\{\mathbf{v}_t^\top \mathbf{V}_{t-1}^{-1} \mathbf{v}_{t-1}, 1\} \leq 2 \log \det \left(\mathbf{I} + \lambda^{-1} \sum_{i=1}^t \mathbf{v}_i \mathbf{v}_i^\top \right).$$

533 *Proof of Lemma 4.8.* First, recall $\bar{\sigma}_{t,k}$ defined in Definition B.2 and the bound of $\bar{\sigma}_{t,k} - \sigma_{t,k}$ provided
 534 in Lemma B.4. If η and m follows the condition presented in Condition 4.1, we have that there exists
 535 a positive constants C_1 such that

$$\begin{aligned} \sum_{i=1}^T \min\{\sigma_{t,a_t}, 1\} &= \sum_{i=1}^T \min\{\bar{\sigma}_{t,a_t}, 1\} + \sum_{i=1}^T (\sigma_{t,a_t} - \bar{\sigma}_{t,a_t}) \\ &\leq \sqrt{T \sum_{i=1}^T \min\{\bar{\sigma}_{t,a_t}^2, 1\}} + C_1 T^{13/6} \sqrt{\log m} m^{-1/6} \lambda^{-2/3} L^{9/2}, \end{aligned}$$

536 where the first term in the inequality on the second line is from Cauchy-Schwartz inequality, and the
 537 second term is from Lemma B.4. From Definition B.2, we have

$$\sum_{i=1}^T \min\{\bar{\sigma}_{t,a_t}^2, 1\} \leq \lambda \sum_{i=1}^T \min\{\mathbf{g}(\mathbf{x}_{t,a_t}, \boldsymbol{\theta}_0)^\top \bar{\mathbf{U}}_{t-1}^{-1} \mathbf{g}(\mathbf{x}_{t,a_t}, \boldsymbol{\theta}_0), 1\} \quad (\text{B.20})$$

$$\leq 2\lambda \log \det \left(\mathbf{I} + \lambda^{-1} \sum_{i=1}^T \mathbf{g}(\mathbf{x}_{t,a_t}; \boldsymbol{\theta}_0) \mathbf{g}(\mathbf{x}_{t,a_t}; \boldsymbol{\theta}_0)^\top / m \right) \quad (\text{B.21})$$

$$= 2\lambda \log \det(\mathbf{I} + \lambda^{-1} \bar{\mathbf{J}}_T \bar{\mathbf{J}}_T^\top / m) \quad (\text{B.22})$$

$$= 2\lambda \log \det(\mathbf{I} + \lambda^{-1} \bar{\mathbf{J}}_T^\top \bar{\mathbf{J}}_T / m) \quad (\text{B.23})$$

$$= 2\lambda \log \det(\mathbf{I} + \lambda^{-1} \mathbf{K}_T) \quad (\text{B.24})$$

538 where (B.20) moves the positive parameter λ outside the min operator and uses the definition of $\bar{\sigma}_{t,k}$ in
 539 Definition B.2, (B.21), (B.22) utilize Lemma B.11, (B.23) is from the fact that $\det(\mathbf{I} + \mathbf{A}\mathbf{A}^\top) = \det(\mathbf{I} + \mathbf{A}^\top \mathbf{A})$, and (B.24) uses the definition of \mathbf{K}_t in
 540 Definition B.2. From Lemma B.7, we have that

$$\log \det(\mathbf{I} + \lambda^{-1} \mathbf{K}_T) \leq \log \det(\mathbf{I} + \lambda^{-1} \mathbf{H}) + 1$$

542 under condition on m and η presented in Theorem 3.5. By taking a union bound we have, with
 543 probability $1 - 2\delta$, that

$$\sum_{i=1}^T \min\{\bar{\sigma}_{t,a_t}, 1\} \leq \sqrt{2\lambda T(\tilde{d} \log(1 + TK) + 1)} + C_1 T^{13/6} \sqrt{\log m} m^{-1/6} \lambda^{-2/3} L^{9/2},$$

544 where we use the definition of \tilde{d} in Definition 3.2. Replacing δ with $\delta/2$ completes the proof. \square

545 C Proof of auxiliary lemmas in Appendix B

546 In this section, we are about to show the proof of the Lemmas used in Appendix B, we will start
 547 with the following NTK Lemmas. Among them, the first is to control the difference between the
 548 parameter learned via Gradient Descent and the theoretical optimal solution to linearized network.

549 **Lemma C.1** (Lemma B.2, Zhou et al. [51]). There exist constants $\{C_i\}_{i=1}^5 > 0$ such that for any
 550 $\delta > 0$, if η, m satisfy that for all $t \in [T]$,

$$\begin{aligned} 2\sqrt{t/\lambda} &\geq C_1 m^{-2} L^{-3/2} [\log(TKL^2/\delta)]^{3/2}, \\ 2\sqrt{t/\lambda} &\leq C_2 \min \left\{ mL^{-6} [\log m]^{-3/2}, (m^2(\lambda\eta)^2 L^{-6} t^{-1} (\log m)^{-1})^{3/8} \right\}, \\ \eta &\leq C_3(m\lambda + tmL)^{-1}, \\ m^{1/6} &\geq C_4 \sqrt{\log m} L^{7/2} t^{7/6} \lambda^{-7/6} (1 + \sqrt{t/\lambda}), \end{aligned}$$

551 then with probability at least $1 - \delta$ over the random initialization of $\boldsymbol{\theta}_0$, for any $t \in [T]$, we have that
 552 $\|\boldsymbol{\theta}_{t-1} - \boldsymbol{\theta}_0\|_2 \leq 2\sqrt{t/(m\lambda)}$ and

$$\begin{aligned} \|\boldsymbol{\theta}_{t-1} - \boldsymbol{\theta}_0 - \bar{\mathbf{U}}_{t-1}^{-1} \bar{\mathbf{J}}_{t-1} \mathbf{r}_{t-1} / m\|_2 \\ \leq (1 - \eta m \lambda)^J \sqrt{t/(m\lambda)} + C_5 m^{-2/3} \sqrt{\log m} L^{7/2} t^{5/3} \lambda^{-5/3} (1 + \sqrt{t/\lambda}). \end{aligned}$$

553 And the next lemma, controls the difference between the function value of neural network and the
554 linearized model:

555 **Lemma C.2** (Lemma 4.1, Cao and Gu [13]). There exist constants $\{C_i\}_{i=1}^3 > 0$ such that for any
556 $\delta > 0$, if τ satisfies that

$$C_1 m^{-3/2} L^{-3/2} [\log(TKL^2/\delta)]^{3/2} \leq \tau \leq C_2 L^{-6} [\log m]^{-3/2},$$

557 then with probability at least $1 - \delta$ over the random initialization of θ_0 , for all $\tilde{\theta}, \hat{\theta}$ satisfying
558 $\|\tilde{\theta} - \theta_0\|_2 \leq \tau, \|\hat{\theta} - \theta_0\|_2 \leq \tau$ and $j \in [TK]$ we have

$$\left| f(\mathbf{x}^j; \tilde{\theta}) - f(\mathbf{x}^j; \hat{\theta}) - \langle \mathbf{g}(\mathbf{x}^j; \hat{\theta}), \tilde{\theta} - \hat{\theta} \rangle \right| \leq C_3 \tau^{4/3} L^3 \sqrt{m \log m}.$$

559 Furthermore, to continue with, next lemma is proposed to control the difference between the gradient
560 and the gradient on the initial point.

561 **Lemma C.3** (Theorem 5, Allen-Zhu et al. [6]). There exist constants $\{C_i\}_{i=1}^3 > 0$ such that for any
562 $\delta \in (0, 1)$, if τ satisfies that

$$C_1 m^{-3/2} L^{-3/2} [\log(TKL^2/\delta)]^{3/2} \leq \tau \leq C_2 L^{-6} [\log m]^{-3/2},$$

563 then with probability at least $1 - \delta$ over the random initialization of θ_0 , for all $\|\theta - \theta_0\|_2 \leq \tau$ and
564 $j \in [TK]$ we have

$$\|\mathbf{g}(\mathbf{x}^j; \theta) - \mathbf{g}(\mathbf{x}^j; \theta_0)\|_2 \leq C_3 \sqrt{\log m} \tau^{1/3} L^3 \|\mathbf{g}(\mathbf{x}^j; \theta_0)\|_2.$$

555 Also, we need the next lemma to control the gradient norm of the neural network with the help of
556 NTK.

557 **Lemma C.4** (Lemma B.3, Cao and Gu [13]). There exist constants $\{C_i\}_{i=1}^3 > 0$ such that for any
558 $\delta > 0$, if τ satisfies that

$$C_1 m^{-3/2} L^{-3/2} [\log(TKL^2/\delta)]^{3/2} \leq \tau \leq C_2 L^{-6} [\log m]^{-3/2},$$

559 then with probability at least $1 - \delta$ over the random initialization of θ_0 , for any $\|\theta - \theta_0\|_2 \leq \tau$ and
560 $j \in [TK]$ we have $\|\mathbf{g}(\mathbf{x}^j; \theta)\|_F \leq C_3 \sqrt{mL}$.

551 Finally, as literally shows, we can also provide bounds on the kernel provided by the linearized model
552 and the NTK kernel if the network is width enough.

553 **Lemma C.5** (Lemma B.1, Zhou et al. [51]). Set $\mathbf{K} = \sum_{t=1}^T \sum_{k=1}^K \mathbf{g}(\mathbf{x}_{t,k}; \theta_0) \mathbf{g}(\mathbf{x}_{t,k}; \theta_0)^T / m$, recall
554 the definition of \mathbf{H} in Definition 3.1 then there exists a constant C_1 such that

$$m \geq C_1 L^6 \log(TKL/\delta) \epsilon^{-4},$$

557 we could get that $\|\mathbf{K} - \mathbf{H}\|_F \leq TK\epsilon$.

556 Equipped with these lemmas, we could continue for our proof.

577 C.1 Proof of Lemma B.4

578 *Proof of Lemma B.4.* Firstly, set $\tau = 2\sqrt{t/(m\lambda)}$, then we have the condition on the network m
579 and learning rate η satisfy all of the condition need from Lemma C.1 to Lemma C.5. Thus from
580 Lemma C.1 we have that there exists $\|\theta_{t-1} - \theta_0\| \leq \tau$, thus from Lemma C.4 we have that there
581 exists positive constant \bar{C}_1 such that $\|\mathbf{g}(\mathbf{x}; \theta_{t-1})\|_2 \leq \bar{C}_1 \sqrt{mL}$, $\|\mathbf{g}(\mathbf{x}; \theta_0)\|_2 \leq \bar{C}_1 \sqrt{mL}$, consider
582 the function defined as

$$\psi(\mathbf{a}, \mathbf{a}_1, \dots, \mathbf{a}_{t-1}) = \sqrt{\mathbf{a}^\top \left(\sum_{i=1}^{t-1} \lambda \mathbf{I} + \mathbf{a}_i \mathbf{a}_i^\top \right)^{-1} \mathbf{a}},$$

583 it is then easy to verify that

$$\begin{aligned} \psi\left(\frac{\mathbf{g}(\mathbf{x}_{t,k}; \theta_{t-1})}{\sqrt{m}}, \frac{\mathbf{g}(\mathbf{x}_{1,a_1}; \theta_0)}{\sqrt{m}}, \dots, \frac{\mathbf{g}(\mathbf{x}_{t-1,a_{t-1}}; \theta_{t-1})}{\sqrt{m}}\right) &= \sigma_{t,k} \\ \psi\left(\frac{\mathbf{g}(\mathbf{x}_{t,k}; \theta_0)}{\sqrt{m}}, \frac{\mathbf{g}(\mathbf{x}_{1,a_1}; \theta_0)}{\sqrt{m}}, \dots, \frac{\mathbf{g}(\mathbf{x}_{t-1,a_{t-1}}; \theta_0)}{\sqrt{m}}\right) &= \bar{\sigma}_{t,k}, \end{aligned}$$

584 then we obtain that the function ψ is defined under the domain $\|\mathbf{a}\|_2 \leq \bar{C}_1\sqrt{L}$, $\|\mathbf{a}_i\|_2 \leq \bar{C}_1\sqrt{L}$ then
585 by taking the derivation w.r.t. ψ^2 , we have that

$$2\psi\partial\psi = (\partial\mathbf{a})^\top \left(\sum_{i=1}^{t-1} \lambda\mathbf{I} + \mathbf{a}_i\mathbf{a}_i^\top \right)^{-1} \mathbf{a} + \mathbf{a}^\top \left(\sum_{i=1}^{t-1} \lambda\mathbf{I} + \mathbf{a}_i\mathbf{a}_i^\top \right)^{-1} \partial\mathbf{a} \\ + \mathbf{a}^\top \left(\sum_{i=1}^{t-1} \lambda\mathbf{I} + \mathbf{a}_i\mathbf{a}_i^\top \right)^{-1} \sum_{i=1}^{t-1} ((\partial\mathbf{a}_i)\mathbf{a}_i^\top + \mathbf{a}_i\partial\mathbf{a}_i^\top) \left(\sum_{i=1}^{t-1} \lambda\mathbf{I} + \mathbf{a}_i\mathbf{a}_i^\top \right)^{-1} \mathbf{a},$$

586 by taking trace with both side and utilizing $\text{tr}(\mathbf{AB}) = \text{tr}(\mathbf{BA})$ and $\text{tr}(\boldsymbol{\alpha}^\top \boldsymbol{\beta}) = \text{tr}(\boldsymbol{\alpha}\boldsymbol{\beta}^\top)$, we have
587 that

$$2\text{tr}(\psi\partial\psi) = \text{tr} \left(2(\partial\mathbf{a})^\top \left(\sum_{i=1}^{t-1} \lambda\mathbf{I} + \mathbf{a}_i\mathbf{a}_i^\top \right)^{-1} \mathbf{a} \right. \\ \left. + 2 \sum_{j=1}^{t-1} (\partial\mathbf{a}_j)^\top \left(\left(\sum_{i=1}^{t-1} \lambda\mathbf{I} + \mathbf{a}_i\mathbf{a}_i^\top \right)^{-1} \mathbf{a} \mathbf{a}^\top \left(\sum_{i=1}^{t-1} \lambda\mathbf{I} + \mathbf{a}_i\mathbf{a}_i^\top \right)^{-1} \mathbf{a}_j \right) \right),$$

588 thus by setting $\mathbf{C} = \left(\sum_{i=1}^{t-1} \lambda\mathbf{I} + \mathbf{a}_i\mathbf{a}_i^\top \right)^{-1}$ for simplicity and decompose $\mathbf{C} = \mathbf{Q}^\top \mathbf{D} \mathbf{Q}$, $\mathbf{b} = \mathbf{Q}\mathbf{a}$
589 where $\mathbf{D} = \text{diag}(\varrho_1, \dots, \varrho_p)$ as the eigen-value of \mathbf{C} , we have that

$$\nabla_{\mathbf{a}}\psi = \frac{\mathbf{Ca}}{\sqrt{\mathbf{a}^\top \mathbf{Ca}}}, \|\nabla_{\mathbf{a}}\psi\|_2 = \sqrt{\frac{\mathbf{a}^\top \mathbf{C}^2 \mathbf{a}}{\mathbf{a}^\top \mathbf{Ca}}} = \sqrt{\frac{\mathbf{b}^\top \mathbf{D}^2 \mathbf{b}}{\mathbf{b}^\top \mathbf{Db}}} = \sqrt{\frac{\sum_{i=1}^d b_i^2 \varrho_i^2}{\sum_{i=1}^d b_i^2 \varrho_i}} \leq 1/\sqrt{\lambda}$$

590 where the last inequality is from the fact that $\mathbf{C} \preceq 1/\lambda\mathbf{I}$, which indicates that all eigen-value $\varrho_i \leq 1/\lambda$,
591 for the same reason, we have

$$\|\nabla_{\mathbf{a}_i}\psi\|_2 = \frac{\|\mathbf{C}\mathbf{a}\mathbf{a}^\top \mathbf{C}\mathbf{a}_j\|_2}{\sqrt{\mathbf{a}^\top \mathbf{C}\mathbf{a}}} \leq \|\mathbf{a}_j\|_2 \frac{\|\mathbf{C}\mathbf{a}\mathbf{a}^\top \mathbf{C}\|_2}{\sqrt{\mathbf{a}^\top \mathbf{C}\mathbf{a}}} = \|\mathbf{a}_j\|_2 \|\mathbf{C}\mathbf{a}\|_2 \frac{\|\mathbf{C}^\top \mathbf{a}\|_2}{\sqrt{\mathbf{a}^\top \mathbf{C}\mathbf{a}}} \leq \|\mathbf{a}_j\| \|\mathbf{a}\| / \sqrt{\lambda}.$$

592 Thus under the domain that $\|\mathbf{a}\|_2 \leq \bar{C}_1\sqrt{L}$, $\|\mathbf{a}_i\|_2 \leq \bar{C}_1\sqrt{L}$, we have that

$$\|\nabla_{\mathbf{a}}\psi\|_2 \leq 1/\sqrt{\lambda}, \|\nabla_{\mathbf{a}_i}\psi\|_2 \leq \bar{C}_1^2 L / \sqrt{\lambda}.$$

593 Then, Lipschitz continuity implies

$$|\sigma_{t,k} - \bar{\sigma}_{t,k}| = \left| \psi \left(\frac{\mathbf{g}(\mathbf{x}_{t,k}; \boldsymbol{\theta}_{t-1})}{\sqrt{m}}, \frac{\mathbf{g}(\mathbf{x}_{1,a_1}; \boldsymbol{\theta}_1)}{\sqrt{m}}, \dots, \frac{\mathbf{g}(\mathbf{x}_{t-1,a_{t-1}}; \boldsymbol{\theta}_{t-1})}{\sqrt{m}} \right) \right. \\ \left. - \psi \left(\frac{\mathbf{g}(\mathbf{x}_{t,k}; \boldsymbol{\theta}_0)}{\sqrt{m}}, \frac{\mathbf{g}(\mathbf{x}_{1,a_1}; \boldsymbol{\theta}_0)}{\sqrt{m}}, \dots, \frac{\mathbf{g}(\mathbf{x}_{t-1,a_{t-1}}; \boldsymbol{\theta}_0)}{\sqrt{m}} \right) \right| \\ \leq \sup\{\|\nabla_{\mathbf{a}}\psi\|_2\} \left\| \frac{\mathbf{g}(\mathbf{x}_{t,k}; \boldsymbol{\theta}_{t-1})}{\sqrt{m}} - \frac{\mathbf{g}(\mathbf{x}_{t,k}; \boldsymbol{\theta}_0)}{\sqrt{m}} \right\|_2 \\ + \sum_{i=1}^{t-1} \sup\{\|\nabla_{\mathbf{a}_i}\psi\|_2\} \left\| \frac{\mathbf{g}(\mathbf{x}_{i,a_i}; \boldsymbol{\theta}_i)}{\sqrt{m}} - \frac{\mathbf{g}(\mathbf{x}_{i,a_i}; \boldsymbol{\theta}_0)}{\sqrt{m}} \right\|_2 \\ \leq \frac{1}{\sqrt{\lambda}} \left\| \frac{\mathbf{g}(\mathbf{x}_{t,k}; \boldsymbol{\theta}_t) - \mathbf{g}(\mathbf{x}_{t,k}; \boldsymbol{\theta}_0)}{\sqrt{m}} \right\|_2 + \frac{\bar{C}_1^2 L}{\sqrt{\lambda}} \sum_{i=1}^{t-1} \left\| \frac{\mathbf{g}(\mathbf{x}_{i,a_i}; \boldsymbol{\theta}_i) - \mathbf{g}(\mathbf{x}_{i,a_i}; \boldsymbol{\theta}_0)}{\sqrt{m}} \right\|_2. \tag{C.1}$$

594 By Lemma C.3 with $\tau = 2\sqrt{t/m\lambda}$, there exist positive constants \bar{C}_2 and \bar{C}_3 so that each gradient
595 difference in (C.1) is bounded by

$$\frac{1}{\sqrt{m}} \|\mathbf{g}(\mathbf{x}; \boldsymbol{\theta}) - \mathbf{g}(\mathbf{x}; \boldsymbol{\theta}_0)\|_2 \leq \bar{C}_2 \sqrt{\log m} \tau^{1/3} L^3 \|\mathbf{g}(\mathbf{x}; \boldsymbol{\theta}_0)\|_2 / \sqrt{m} \\ \leq \bar{C}_3 \sqrt{\log m}^{1/6} m^{-1/6} \lambda^{-1/6} L^{7/2}.$$

596 Thus, since we obtain that there exists constant C_5 such that

$$|\sigma_{t,k} - \bar{\sigma}_{t,k}| \leq C_1 \sqrt{\log m} t^{7/6} m^{-1/6} \lambda^{-2/3} L^{9/2},$$

597 where we use the fact that $C_1 = \max\{\bar{C}_3, \bar{C}_3 \bar{C}_1^2\}$ and $L \geq 1$ to merge the first term into the
598 summation. This inequality is based on Lemma C.1, Lemma C.3 and Lemma C.4, thus it holds with
599 probability at least $1 - 3\delta$. Replacing δ with $\delta/3$ completes the proof. \square

600 C.2 Proof of Lemma B.5

601 *Proof of Lemma B.5.* Setting $\tau = 2\sqrt{t/m\lambda}$, we have the condition on the network m and learning
602 rate η satisfy all of the condition needed by Lemmas C.1 to C.5. From Lemma C.1 we have
603 $\theta_{t-1} - \theta_0 \leq \tau$. Then, by Lemma C.2, there exists a constant C_1 such that

$$|f(\mathbf{x}_{t,k}; \theta_{t-1}) - \langle \mathbf{g}(\mathbf{x}_{t,k}; \theta_0), \theta_{t-1} - \theta_0 \rangle| \leq C_1 t^{2/3} m^{-1/6} \lambda^{-2/3} L^3 \sqrt{\log m}, \quad (\text{C.2})$$

604 Using the bound on $\theta_{t-1} - \theta_0 - \bar{\mathbf{U}}_{t-1}^{-1} \bar{\mathbf{J}}_{t-1} \mathbf{r}_{t-1}/m$ provided in Lemma C.1 and the norm of gradient
605 bound given in Lemma C.4, we have that there exist positive constants C_1, \bar{C}_2 such that

$$\begin{aligned} & |\langle \mathbf{g}(\mathbf{x}_{t,k}; \theta_0), \theta_{t-1} - \theta_0 \rangle - \langle \mathbf{g}(\mathbf{x}_{t,k}; \theta_0), \bar{\mathbf{U}}_{t-1}^{-1} \bar{\mathbf{J}}_{t-1} \mathbf{r}_{t-1}/m \rangle| \\ & \leq \|\mathbf{g}(\mathbf{x}_{t,k})\|_2 \|\theta_{t-1} - \theta_0 - \bar{\mathbf{U}}_{t-1}^{-1} \bar{\mathbf{J}}_{t-1} \mathbf{r}_{t-1}/m\|_2 \\ & \leq \bar{C}_1 \sqrt{mL} \left((1 - \eta m \lambda)^J \sqrt{t/(m\lambda)} + \bar{C}_2 m^{-2/3} \sqrt{\log m} L^{7/2} t^{5/3} \lambda^{-5/3} (1 + \sqrt{t/\lambda}) \right) \\ & = C_2 (1 - \eta m \lambda)^J \sqrt{tL/\lambda} + C_3 m^{-1/6} \sqrt{\log m} L^4 t^{5/3} \lambda^{-5/3} (1 + \sqrt{t/\lambda}), \end{aligned} \quad (\text{C.3})$$

606 where $C_2 = \bar{C}_1, C_3 = \bar{C}_1 \bar{C}_2$. Combining (C.2) and (C.3), we have

$$\begin{aligned} |f(\mathbf{x}_{t,k}) - \langle \mathbf{g}(\mathbf{x}_{t,k}; \theta_0), \bar{\mathbf{U}}_{t-1}^{-1} \bar{\mathbf{J}}_{t-1} \mathbf{r}_{t-1}/m \rangle| & \leq C_1 t^{2/3} m^{-1/6} \lambda^{-2/3} L^3 \sqrt{\log m} \\ & + C_2 (1 - \eta m \lambda)^J \sqrt{tL/\lambda} \\ & + C_3 m^{-1/6} \sqrt{\log m} L^4 t^{5/3} \lambda^{-5/3} (1 + \sqrt{t/\lambda}), \end{aligned}$$

607 which holds with probability $1 - 3\delta$ with a union bound (Lemma C.4, Lemma C.1, and Lemma C.2).
608 Replacing δ with $\delta/3$ completes the proof. \square

609 C.3 Proof of Lemma B.7

610 *Proof of Lemma B.7.* From the definition of \mathbf{K}_t , we have that

$$\begin{aligned} \log \det(\mathbf{I} + \lambda^{-1} \mathbf{K}_t) & = \log \det \left(\mathbf{I} + \sum_{i=1}^t \mathbf{g}(\mathbf{x}_{i,a_i}; \theta_0) \mathbf{g}(\mathbf{x}_{i,a_i}; \theta_0)^\top / (m\lambda) \right) \\ & \leq \log \det \left(\mathbf{I} + \sum_{t=1}^T \sum_{k=1}^K \mathbf{g}(\mathbf{x}_{i,a_i}; \theta_0) \mathbf{g}(\mathbf{x}_{i,a_i}; \theta_0)^\top / (m\lambda) \right) \end{aligned} \quad (\text{C.4})$$

$$= \log \det(\mathbf{I} + \mathbf{K}/\lambda) \quad (\text{C.5})$$

$$\leq \log \det(\mathbf{I} + \mathbf{H}/\lambda + (\mathbf{H} - \mathbf{K})\lambda) + T(\lambda - 1) \quad (\text{C.6})$$

$$\leq \log \det(\mathbf{I} + \mathbf{H}/\lambda) + \langle (\mathbf{I} + \mathbf{H}/\lambda)^{-1}, (\mathbf{K} - \mathbf{H})/\lambda \rangle \quad (\text{C.7})$$

$$\leq \log \det(\mathbf{I} + \mathbf{H}/\lambda) + \|(\mathbf{I} + \mathbf{H}/\lambda)^{-1}\|_F \|(\mathbf{K} - \mathbf{H})\|_F / \lambda \quad (\text{C.8})$$

$$\leq \log \det(\mathbf{I} + \mathbf{H}/\lambda) + \sqrt{TK} \|(\mathbf{K} - \mathbf{H})\|_F \quad (\text{C.9})$$

611 where the inequality (C.4) is because the double summation on the second line contains more
612 elements than the summation on the first line. The inequality (C.5) utilize the definition of \mathbf{K} in
613 Lemma C.5 and \mathbf{H} in Definition B.1, inequality (C.6) is from the convexity of $\log \det(\cdot)$ function, and
614 inequality (C.7) is from the fact that $\langle \mathbf{A}, \mathbf{B} \rangle \leq \|\mathbf{A}\|_F \|\mathbf{B}\|_F$. Inequality (C.8) is from the fact that
615 $\|\mathbf{A}\|_F \leq \sqrt{TK} \|\mathbf{A}\|_2$ if $\mathbf{A} \in \mathbb{R}^{TK \times TK}$ and $\lambda \geq 0$, inequality (C.9) utilizes Lemma C.5 by setting
616 $\epsilon = (TK)^{-3/2}$ with $m \geq C_1 L^6 T^6 K^6 \log(TKL/\delta)$, where we conclude our proof. \square

617 **C.4 Proof of Lemma B.9**

618 *Proof of Lemma B.9.* Set τ in Lemma C.4 as $2\sqrt{t/(m\lambda)}$. Then the network width m and learning
 619 rate η satisfy all of the condition needed by Lemma C.1 to C.5. Hence, there exists C_1 such that
 620 $\|\mathbf{g}(\mathbf{x}; \boldsymbol{\theta})\|_2 \leq \|\mathbf{g}(\mathbf{x}; \boldsymbol{\theta})\|_F \leq C_1\sqrt{mL}$ for all \mathbf{x} , since it is easy to verify that $\mathbf{U}_t^{-1} \preceq \lambda^{-1}\mathbf{I}$. Thus we
 621 have that for all $t \in [T], k \in [K]$,

$$\sigma_{t,k}^2 = \lambda \mathbf{g}^\top(\mathbf{x}_{t,k}; \boldsymbol{\theta}_{t-1}) \mathbf{U}_{t-1}^{-1} \mathbf{g}(\mathbf{x}_{t,k}; \boldsymbol{\theta}_{t-1}) / m \leq \|\mathbf{g}(\mathbf{x}_{t,k}; \boldsymbol{\theta}_{t-1})\|_2^2 / m \leq C_5^2 L.$$

622 Therefore, we could get that $\sigma_{t,k} \leq C_1\sqrt{L}$, with probability $1 - 2\delta$ by taking a union bound
 623 (Lemmas C.1 and C.4). Replacing δ with $\delta/2$ completes the proof. \square