

---

# Neural Contextual Bandits with Deep Representation and Shallow Exploration

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 We study neural contextual bandits, a general class of contextual bandits, where  
2 each context-action pair is associated with a raw feature vector, but the specific  
3 reward generating function is unknown. We propose a novel learning algorithm  
4 that transforms the raw feature vector using the last hidden layer of a deep ReLU  
5 neural network (deep representation learning), and uses an upper confidence bound  
6 (UCB) approach to explore in the last linear layer (shallow exploration). We prove  
7 that under standard assumptions, our proposed algorithm achieves  $\tilde{O}(\sqrt{T})$  finite-  
8 time regret, where  $T$  is the learning time horizon. Compared with existing neural  
9 contextual bandit algorithms, our approach is computationally much more efficient  
10 since it only needs to explore in the last layer of the deep neural network.

## 11 1 Introduction

12 Multi-armed bandits (MAB) [9, 8, 30] are a class of online decision-making problems where an  
13 agent needs to learn to maximize its expected cumulative reward while repeatedly interacting with a  
14 partially known environment. Based on a bandit algorithm (also called a strategy or policy), in each  
15 round, the agent adaptively chooses an arm, and then observes and receives a reward associated with  
16 that arm. Since only the reward of the chosen arm will be observed (bandit information feedback),  
17 a good bandit algorithm has to deal with the exploration-exploitation dilemma: trade-off between  
18 pulling the best arm based on existing knowledge/history data (exploitation) and trying the arms that  
19 have not been fully explored (exploration).

20 In many real-world applications, the agent will also be able to access detailed contexts associated  
21 with the arms. For example, when a company wants to choose an advertisement to present to a user,  
22 the recommendation will be much more accurate if the company takes into consideration the contents,  
23 specifications, and other features of the advertisements in the arm set as well as the profile of the user.  
24 To encode the contextual information, contextual bandit models and algorithms have been developed,  
25 and widely studied both in theory and in practice [19, 39, 34, 16, 1]. Most existing contextual bandit  
26 algorithms assume that the expected reward of an arm at a context is a linear function in a known  
27 context-action feature vector, which leads to many useful algorithms such as LinUCB [16], OFUL [1],  
28 etc. The representation power of the linear model can be limited in applications such as marketing,  
29 social networking, clinical studies, etc., where the rewards are usually counts or binary variables. The  
30 linear contextual bandit problem has also been extended to richer classes of parametric bandits such  
31 as the generalized linear bandits [24, 35] and kernelised bandits [44, 15].

32 With the prevalence of deep neural networks (DNNs) and their phenomenal performances in many  
33 machine learning tasks [32, 25], there has emerged a line of work that employs DNNs to increase the  
34 representation power of contextual bandit algorithms [5, 38, 17, 49, 52, 20, 51]. The problems they  
35 solve are usually referred to as *neural contextual bandits*. For example, Zhou et al. [52] developed  
36 the NeuralUCB algorithm, which can be viewed as a natural extension of LinUCB [16, 1], where they

37 use the output of a deep neural network with the feature vector as input to approximate the reward.  
 38 Zhang et al. [51] adapted neural networks in Thompson Sampling [43, 14, 40] for both exploration  
 39 and exploitation and proposed NeuralTS. For a fixed time horizon  $T$ , it has been proved that both  
 40 NeuralUCB and NeuralTS achieve a  $O(\tilde{d}\sqrt{T})$  regret bound, where  $\tilde{d}$  is the effective dimension of a  
 41 neural tangent kernel matrix which can potentially scale with  $O(TK)$  for  $K$ -armed bandits. This  
 42 high complexity is mainly due to that the exploration is performed over the entire huge neural network  
 43 parameter space, which is inefficient and even infeasible when the number of neurons is large. A more  
 44 realistic and efficient way of learning neural contextual bandits may be to just explore different arms  
 45 using the last layer as the exploration parameter. More specifically, Riquelme et al. [38] provided  
 46 an extensive empirical study of benchmark algorithms for contextual-bandits through the lens of  
 47 Thompson Sampling, which suggests decoupling representation learning and uncertainty estimation  
 48 improves performance.

49 In this paper, we show that the decoupling of representation learning and the exploration can be  
 50 theoretically validated. We study a new neural contextual bandit algorithm, which learns a mapping  
 51 to transform the raw features associated with each context-action pair using a deep neural network  
 52 (*deep representation*), and then performs an upper confidence bound (UCB)-type exploration over the  
 53 linear output layer of the network (*shallow exploration*). We prove a sublinear regret of the proposed  
 54 algorithm by exploiting the UCB exploration techniques in linear contextual bandits [1] and the  
 55 analysis of deep overparameterized neural networks using neural tangent kernels [27]. Our theory  
 56 confirms the empirically observed effectiveness of decoupling the deep representation learning and  
 57 the UCB exploration in contextual bandits [38, 49].

58 **Contributions** we summarize the main contributions of this paper as follows.

- 59 • We propose a contextual bandit algorithm, Neural-LinUCB, for solving a general class of con-  
 60 textual bandit problems without knowing the specific reward generating function. The proposed  
 61 algorithm learns a deep representation to transform the raw feature vectors and performs UCB-type  
 62 exploration in the last layer of the neural network, which we refer to as deep representation and  
 63 shallow exploration. Compared with LinUCB [34, 16] and neural bandits such as NeuralUCB [52]  
 64 and NeuralTS [51], our algorithm enjoys the best of two worlds: strong expressiveness due to the  
 65 deep representation and computational efficiency due to the shallow exploration.
- 66 • Despite the usage of a DNN as the feature mapping, we prove a  $\tilde{O}(\sqrt{T})$  regret for the proposed  
 67 Neural-LinUCB algorithm, which matches the regret bound of linear contextual bandits [16, 1].  
 68 To the best of our knowledge, this is the first work that theoretically shows the convergence of  
 69 bandits algorithms under the scheme of deep representation and shallow exploration. It is notable  
 70 that a similar scheme called Neural-Linear was proposed by Riquelme et al. [38] for Thompson  
 71 sampling algorithms, and they empirically showed that decoupling representation learning and  
 72 uncertainty estimation improves the performance. Our work confirms this observation from a  
 73 theoretical perspective.
- 74 • We conduct experiments on contextual bandit problems based on real-world datasets, demon-  
 75 strating a better performance and computational efficiency of Neural-LinUCB over LinUCB and  
 76 NeuralUCB, which well aligns with our theory.

## 77 1.1 Additional related work

78 There is a line of related work to ours on the recent advance in the optimization and generalization  
 79 analysis of deep neural networks. In particular, Jacot et al. [27] first introduced the neural tangent  
 80 kernel (NTK) to characterize the training dynamics of network outputs in the infinite width limit.  
 81 From the notion of NTK, a fruitful line of research emerged and showed that loss functions of  
 82 deep neural networks trained by (stochastic) gradient descent can converge to the global minimum  
 83 [22, 4, 21, 54, 53]. The generalization bounds for overparameterized deep neural networks are also  
 84 established in Arora et al. [6, 7], Allen-Zhu et al. [3], Cao and Gu [12, 13]. Recently, the NTK based  
 85 analysis is also extended to the study of sequential decision problems including bandits [52, 51], and  
 86 reinforcement learning algorithms [11, 36, 45, 47].

87 Our algorithm is also different from Langford and Zhang [29], Agarwal et al. [2] which reduce the  
 88 bandit problem to supervised learning. Moreover, their algorithms need to access an oracle that  
 89 returns the optimal policy in a policy class given a sequence of context and reward vectors, whose  
 90 regret depends on the VC-dimension of the policy class.

91 **Notation** We use  $[k]$  to denote a set  $\{1, \dots, k\}$ ,  $k \in \mathbb{N}^+$ .  $\|\mathbf{x}\|_2 = \sqrt{\mathbf{x}^\top \mathbf{x}}$  is the Euclidean norm of  
92 a vector  $\mathbf{x} \in \mathbb{R}^d$ . For a matrix  $\mathbf{W} \in \mathbb{R}^{m \times n}$ , we denote by  $\|\mathbf{W}\|_2$  and  $\|\mathbf{W}\|_F$  its operator norm  
93 and Frobenius norm respectively. For a semi-definite matrix  $\mathbf{A} \in \mathbb{R}^{d \times d}$  and a vector  $\mathbf{x} \in \mathbb{R}^d$ , we  
94 denote the Mahalanobis norm as  $\|\mathbf{x}\|_{\mathbf{A}} = \sqrt{\mathbf{x}^\top \mathbf{A} \mathbf{x}}$ . Throughout this paper, we reserve the notations  
95  $\{C_i\}_{i=0,1,\dots}$  to represent absolute positive constants that are independent of problem parameters such  
96 as dimension, sample size, iteration number, step size, network length and so on. The specific values  
97 of  $\{C_i\}_{i=0,1,\dots}$  can be different in different context. For a parameter of interest  $T$  and a function  
98  $f(T)$ , we use notations such as  $O(f(T))$  and  $\Omega(f(T))$  to hide constant factors and  $\tilde{O}(f(T))$  to hide  
99 constant and logarithmic dependence of  $T$ .

## 100 2 Preliminaries

101 In this section, we provide the background of contextual bandits and deep neural networks.

### 102 2.1 Linear contextual bandits

103 A contextual bandit is characterized by a tuple  $(\mathcal{S}, \mathcal{A}, r)$ , where  $\mathcal{S}$  is the context (state) space,  $\mathcal{A}$  is the  
104 arm (action) space, and  $r$  encodes the unknown *reward generating function* at all context-arm pairs.  
105 A learning agent, who knows  $\mathcal{S}$  and  $\mathcal{A}$  but does not know the true reward  $r$  (values bounded in  $(0, 1)$   
106 for simplicity), needs to interact with the contextual bandit for  $T$  rounds. At each round  $t = 1, \dots, T$ ,  
107 the agent first observes a context  $s_t \in \mathcal{S}$  chosen by the environment; then it needs to adaptively select  
108 an arm  $a_t \in \mathcal{A}$  based on its past observations; finally it receives a reward  $\hat{r}_t(\mathbf{x}_{s_t, a_t}) = r(\mathbf{x}_{s_t, a_t}) + \xi_t$ ,  
109 where  $\mathbf{x}_{s_t, a_t} \in \mathbb{R}^d$  is a known feature vector for context-arm pair  $(s_t, a_t) \in \mathcal{S} \times \mathcal{A}$ , and  $\xi_t$  is a random  
110 noise with zero mean. The agent’s objective is to maximize its expected total reward over these  $T$   
111 rounds, which is equivalent to minimizing the pseudo regret [8]:

$$R_T = \mathbb{E} \left[ \sum_{t=1}^T (\hat{r}(\mathbf{x}_{s_t, a_t^*}) - \hat{r}(\mathbf{x}_{s_t, a_t})) \right], \quad (2.1)$$

112 where  $a_t^* \in \operatorname{argmax}_{a \in \mathcal{A}} \{r(\mathbf{x}_{s_t, a})\} = \mathbb{E}[\hat{r}(\mathbf{x}_{s_t, a})]$ . To simplify the exposition, we use  $\mathbf{x}_{t,a}$  to denote  
113  $\mathbf{x}_{s_t, a}$  since it only depends on the round index  $t$  in most bandit problems, and we assume  $\mathcal{A} = [K]$ .

114 In some practical problems, the agent has a prior knowledge that the reward-generating function  
115  $r$  has some specific parametric form. For instance, in linear contextual bandits, the agent knows  
116 that  $r(\mathbf{x}_{s_t, a_t}) = \mathbf{x}_{s_t, a_t}^\top \boldsymbol{\theta}^*$  for some unknown weight vector  $\boldsymbol{\theta}^* \in \mathbb{R}^d$ . One provably sample efficient  
117 algorithm for linear contextual bandits is Linear Upper Confidence Bound (LinUCB) [1]. Specifically,  
118 at each round  $t$ , LinUCB chooses action by the following strategy

$$a_t = \operatorname{argmax}_{a \in [K]} \left\{ \mathbf{x}_{t,a}^\top \boldsymbol{\theta}_t + \alpha_t \|\mathbf{x}_{t,a}\|_{\mathbf{A}_t^{-1}} \right\},$$

119 where  $\boldsymbol{\theta}_t$  is a point estimate of  $\boldsymbol{\theta}^*$ ,  $\mathbf{A}_t = \lambda \mathbf{I} + \sum_{i=1}^t \mathbf{x}_{i,a_i} \mathbf{x}_{i,a_i}^\top$  with some  $\lambda > 0$  is a matrix  
120 defined based on the historical context-arm pairs, and  $\alpha_t > 0$  is a tuning parameter that controls the  
121 exploration rate in LinUCB.

### 122 2.2 Deep neural networks

123 In this paper, we use  $f(\mathbf{x})$  to denote a neural network with input data  $\mathbf{x} \in \mathbb{R}^d$ . Let  $L$  be the number  
124 of hidden layers and  $\mathbf{W}_l \in \mathbb{R}^{m_l \times m_{l-1}}$  be the weight matrices in the  $l$ -th layer, where  $l = 1, \dots, L$ ,  
125  $m_1 = \dots = m_{L-1} = m$  and  $m_0 = m_L = d$ . Then a  $L$ -hidden layer neural network is defined as

$$f(\mathbf{x}) = \sqrt{m} \boldsymbol{\theta}^{*\top} \sigma_L(\mathbf{W}_L \sigma_{L-1}(\mathbf{W}_{L-1} \cdots \sigma_1(\mathbf{W}_1 \mathbf{x}) \cdots)), \quad (2.2)$$

126 where  $\sigma_l$  is an activation function and  $\boldsymbol{\theta}^* \in \mathbb{R}^d$  is the weight of the output layer. To simplify the  
127 presentation, we will assume  $\sigma_1 = \sigma_2 = \dots = \sigma_L = \sigma$  is the ReLU activation function, i.e.,  
128  $\sigma(x) = \max\{0, x\}$  for  $x \in \mathbb{R}$ . We denote  $\mathbf{w} = (\operatorname{vec}(\mathbf{W}_1)^\top, \dots, \operatorname{vec}(\mathbf{W}_L)^\top)^\top$ , which is the  
129 concatenation of the vectorized weight parameters of all hidden layers of the neural network. We also  
130 write  $f(\mathbf{x}; \boldsymbol{\theta}^*, \mathbf{w}) = f(\mathbf{x})$  in order to explicitly specify the weight parameters of neural network  $f$ . It

131 is easy to show that the dimension  $p$  of vector  $\mathbf{w}$  satisfies  $p = (L - 2)m^2 + 2md$ . To simplify the  
 132 notation, we define  $\phi(\mathbf{x}; \mathbf{w})$  as the output of the  $L$ -th hidden layer of neural network  $f$ .

$$\phi(\mathbf{x}; \mathbf{w}) = \sqrt{m}\sigma(\mathbf{W}_L\sigma(\mathbf{W}_{L-1}\cdots\sigma(\mathbf{W}_1\mathbf{x})\cdots)). \quad (2.3)$$

133 Note that  $\phi(\mathbf{x}; \mathbf{w})$  itself can also be viewed as a neural network with vector-valued outputs.

### 134 3 Deep Representation and Shallow Exploration

135 The linear parametric form in linear contextual bandits might produce biased estimates of the reward  
 136 due to the lack of representation power [42, 38]. In contrast, it is well known that deep neural networks  
 137 are powerful enough to approximate an arbitrary function [18]. Therefore, a natural extension of  
 138 linear contextual bandits is to use a deep neural network to approximate the reward generating  
 139 function  $r(\cdot)$ . Nonetheless, DNNs usually have a prohibitively large dimension for weight parameters,  
 140 which makes the exploration in neural networks based UCB algorithm inefficient [28, 52].

141 In this work, we study a neural contextual bandit algorithm, where the hidden layers of a deep neural  
 142 network are used to represent the features and the exploration is only performed in the last layer of the  
 143 neural network. In particular, we assume that the reward generating function  $r(\cdot)$  can be expressed as  
 144 the inner product between a deep represented feature vector and an exploration weight parameter,  
 145 namely,  $r(\cdot) = \langle \boldsymbol{\theta}^*, \boldsymbol{\psi}(\cdot) \rangle$ , where  $\boldsymbol{\theta}^* \in \mathbb{R}^d$  is some weight parameter and  $\boldsymbol{\psi}(\cdot)$  is an unknown feature  
 146 mapping. This decoupling of the representation and the exploration will achieve the best of both  
 147 worlds: efficient exploration in shallow (linear) models and high expressive power of deep models.  
 148 To learn the unknown feature mapping, we propose to use a neural network to approximate it. In  
 149 what follows, we will describe a neural contextual bandit algorithm that uses the output of the last  
 150 hidden layer of a neural network to transform the raw feature vectors (*deep representation*) and  
 151 performs UCB-type exploration in the last layer of the neural network (*shallow exploration*). Since  
 152 the exploration is performed only in the last linear layer, we call this procedure Neural-LinUCB,  
 153 which is displayed in Algorithm 1.

154 Specifically, in round  $t$ , the agent receives an action set with raw features  $\mathcal{X}_t = \{\mathbf{x}_{t,1}, \dots, \mathbf{x}_{t,K}\}$ .  
 155 Then the agent chooses an arm  $a_t$  that maximizes the following upper confidence bound:

$$a_t = \operatorname{argmax}_{k \in [K]} \left\{ \langle \phi(\mathbf{x}_{t,k}; \mathbf{w}_{t-1}), \boldsymbol{\theta}_{t-1} \rangle + \alpha_t \|\phi(\mathbf{x}_{t,k}; \mathbf{w}_{t-1})\|_{\mathbf{A}_{t-1}^{-1}} \right\}, \quad (3.1)$$

156 where  $\boldsymbol{\theta}_{t-1}$  is a point estimate of the unknown weight in the last layer,  $\phi(\mathbf{x}; \mathbf{w})$  is defined as in (2.3),  
 157  $\mathbf{w}_{t-1}$  is an estimate of all the weight parameters in the hidden layers of the neural network,  $\alpha_t > 0$  is  
 158 the algorithmic parameter controlling the exploration, and  $\mathbf{A}_t$  is a matrix defined based on historical  
 159 transformed features:

$$\mathbf{A}_t = \lambda \mathbf{I} + \sum_{i=1}^t \phi(\mathbf{x}_{i,a_i}; \mathbf{w}_{i-1}) \phi(\mathbf{x}_{i,a_i}; \mathbf{w}_{i-1})^\top, \quad (3.2)$$

160 and  $\lambda > 0$ . After pulling arm  $a_t$ , the agent will observe a noisy reward  $\hat{r}_t := \hat{r}(\mathbf{x}_{t,a_t})$  defined as

$$\hat{r}(\mathbf{x}_{t,k}) = r(\mathbf{x}_{t,k}) + \xi_t, \quad (3.3)$$

161 where  $\xi_t$  is an independent  $\nu$ -subGaussian random noise for some  $\nu > 0$  and  $r(\cdot)$  is an unknown  
 162 reward function. In this paper, we will interchangeably use notation  $\hat{r}_t$  to denote the reward received  
 163 at the  $t$ -th step and an equivalent notation  $\hat{r}(\mathbf{x})$  to express its dependence on the feature vector  $\mathbf{x}$ .

164 Upon receiving the reward  $\hat{r}_t$ , the agent updates its estimate  $\boldsymbol{\theta}_t$  of the output layer weight by using  
 165 the same  $\ell^2$ -regularized least-squares estimate in linear contextual bandits [1]. In particular, we have

$$\boldsymbol{\theta}_t = \mathbf{A}_t^{-1} \mathbf{b}_t, \quad (3.4)$$

166 where  $\mathbf{b}_t = \sum_{i=1}^t \hat{r}_i \phi(\mathbf{x}_{i,a_i}; \mathbf{w}_{i-1})$ .

167 To save the computation, the neural network  $\phi(\cdot; \mathbf{w}_t)$  will be updated once every  $H$  steps. Therefore,  
 168 we have  $\mathbf{w}_{(q-1)H+1} = \dots = \mathbf{w}_{qH}$  for  $q = 1, 2, \dots$ . We call the time steps  $\{(q-1)H+1, \dots, qH\}$   
 169 an epoch with length  $H$ . At time step  $t = Hq$ , for any  $q = 1, 2, \dots$ , Algorithm 1 will retrain the

170 neural network based on all the historical data. In Algorithm 2, our goal is to minimize the following  
 171 empirical loss function:

$$\mathcal{L}_q(\mathbf{w}) = \sum_{i=1}^{qH} (\boldsymbol{\theta}_i^\top \phi(\mathbf{x}_{i,a_i}; \mathbf{w}) - \widehat{r}_i)^2. \quad (3.5)$$

172 In practice, one can further save computational cost by only feeding data  $\{\mathbf{x}_{i,a_i}, \widehat{r}_i, \boldsymbol{\theta}_i\}_{i=(q-1)H+1}^{qH}$   
 173 from the  $q$ -th epoch into Algorithm 2 to update the parameter  $\mathbf{w}_t$ , which does not hurt the performance  
 174 since the historical information has been encoded into the estimate of  $\boldsymbol{\theta}_i$ . In this paper, we will  
 175 perform the following gradient descent step

$$\mathbf{w}_q^{(s)} = \mathbf{w}_q^{(s-1)} - \eta_q \nabla_{\mathbf{w}} \mathcal{L}_q(\mathbf{w}_q^{(s-1)}).$$

176 for  $s = 1, \dots, n$ , where  $\mathbf{w}_q^{(0)} = \mathbf{w}^{(0)}$  is chosen as the same random initialization point. We will  
 177 discuss more about the initial point  $\mathbf{w}^{(0)}$  in the next paragraph. Then Algorithm 2 outputs  $\mathbf{w}_q^{(n)}$  and  
 178 we set it as the updated weight parameter  $\mathbf{w}_{Hq+1}$  in Algorithm 1. In the next round, the agent will  
 179 receive another action set  $\mathcal{X}_{t+1}$  with raw feature vectors and repeat the above steps to choose the  
 180 sub-optimal arm and update estimation for contextual parameters.

181 **Initialization:** Recall that  $\mathbf{w}$  is the collection of all hidden layer weight parameters of the neural  
 182 network. We will follow the same initialization scheme as used in Zhou et al. [52], where each entry  
 183 of the weight matrices follows some Gaussian distribution. Specifically, for any  $l \in \{1, \dots, L-1\}$ ,

184 we set  $\mathbf{W}_l = \begin{bmatrix} \mathbf{W} & \mathbf{0} \\ \mathbf{0} & \mathbf{W} \end{bmatrix}$ , where each entry of  $\mathbf{W}$  follows distribution  $N(0, 4/m)$  independently; for

185  $\mathbf{W}_L$ , we set it as  $\begin{bmatrix} \mathbf{V} & \\ & -\mathbf{V} \end{bmatrix}$ , where each entry of  $\mathbf{V}$  follows distribution  $N(0, 2/m)$  independently.

186 **Comparison with LinUCB and NeuralUCB:** Compared with linear contextual bandits in Sec-  
 187 tion 2.1, Algorithm 1 has a distinct feature that it learns a deep neural network to obtain a deep  
 188 representation of the raw data vectors and then performs UCB exploration. This deep representa-  
 189 tion allows our algorithm to characterize more intrinsic and latent information about the raw data  
 190  $\{\mathbf{x}_{t,k}\}_{t \in [T], k \in [K]} \subset \mathbb{R}^d$ . However, the increased complexity of the feature mapping  $\phi(\cdot; \mathbf{w})$  also  
 191 introduces great hardness in training. For instance, a recent work by Zhou et al. [52] also stud-  
 192 ied the neural contextual bandit problem, but different from (3.1), their algorithm (NeuralUCB)  
 193 performs the UCB exploration on the entire network parameter space, which is  $\mathbb{R}^{\widetilde{p}+d}$ , where  
 194  $\widetilde{p} = m + md + (L-1)m^2$ . Note that in Zhou et al. [52], they need to compute the inverse  
 195 of a matrix  $\mathbf{Z}_t \in \mathbb{R}^{(\widetilde{p}+d) \times (\widetilde{p}+d)}$ , which is defined in a similar way to the matrix  $\mathbf{A}_t$  in our paper  
 196 except that  $\mathbf{Z}_t$  is defined based on the gradient of the network instead of the output of the last hidden  
 197 layer as in (3.2). In sharp contrast,  $\mathbf{A}_t$  in our paper is only of size  $d \times d$  and thus is much more  
 198 efficient and practical in implementation, which will be seen from our experiments in later sections.

199 We note that there is also a similar algorithm to our Neural-LinUCB presented in Deshmukh et al.  
 200 [20], where they studied the self-supervised learning loss in contextual bandits with neural network  
 201 representation for computer vision problems. However, no regret analysis has been provided. When  
 202 the feature mapping  $\phi(\cdot; \mathbf{w})$  is an identity function, the problem reduces to linear contextual bandits  
 203 where we directly use  $\mathbf{x}_t$  as the feature vector. In this case, it is easy to see that Algorithm 1 reduces  
 204 to LinUCB [16] since we do not need to learn the representation parameter  $\mathbf{w}$  anymore.

205 **Comparison with Neural-Linear:** The high-level idea of decoupling the representation and explo-  
 206 ration in our algorithm is also similar to that of the Neural-Linear algorithm [38, 49], which trains a  
 207 deep neural network to learn a representation of the raw feature vectors, and then uses a Bayesian  
 208 linear regression to estimate the uncertainty in the bandit problem. However, these two algorithms  
 209 are significantly different since Neural-Linear [38] is a Thompson sampling based algorithm that  
 210 uses posterior sampling to estimate the weight parameter  $\boldsymbol{\theta}^*$  via Bayesian linear regression, whereas  
 211 Neural-LinUCB adopts upper confidence bound based techniques to estimate the weight  $\boldsymbol{\theta}^*$ . Never-  
 212 theless, both algorithms share the same idea of deep representation and shallow exploration, and we  
 213 view our Neural-LinUCB algorithm as one instantiation of the Neural-Linear scheme.

## 214 4 Main Results

215 To analyze the regret bound of Algorithm 1, we first lay down some important assumptions on the  
 216 neural contextual bandit model.

---

**Algorithm 1** Deep Representation and Shallow Exploration (Neural-LinUCB)

---

- 1: **Input:** regularization parameter  $\lambda > 0$ , number of total steps  $T$ , episode length  $H$ , exploration parameters  $\{\alpha_t > 0\}_{t \in [T]}$
  - 2: **Initialization:**  $\mathbf{A}_0 = \lambda \mathbf{I}$ ,  $\mathbf{b}_0 = \mathbf{0}$ ; entries of  $\boldsymbol{\theta}_0$  follow  $N(0, 1/d)$ , and  $\mathbf{w}^{(0)}$  is initialized as described in Section 3;  $q = 1$ ;  $\mathbf{w}_0 = \mathbf{w}^{(0)}$
  - 3: **for**  $t = 1, \dots, T$  **do**
  - 4:   receive feature vectors  $\{\mathbf{x}_{t,1}, \dots, \mathbf{x}_{t,K}\}$
  - 5:   choose arm  $a_t = \operatorname{argmax}_{k \in [K]} \boldsymbol{\theta}_{t-1}^\top \boldsymbol{\phi}(\mathbf{x}_{t,k}; \mathbf{w}_{t-1}) + \alpha_t \|\boldsymbol{\phi}(\mathbf{x}_{t,k}; \mathbf{w}_{t-1})\|_{\mathbf{A}_{t-1}^{-1}}$ , and obtain reward  $\hat{r}_t$
  - 6:   update  $\mathbf{A}_t$  and  $\mathbf{b}_t$  as follows:  
     $\mathbf{A}_t = \mathbf{A}_{t-1} + \boldsymbol{\phi}(\mathbf{x}_{t,a_t}; \mathbf{w}_{t-1}) \boldsymbol{\phi}(\mathbf{x}_{t,a_t}; \mathbf{w}_{t-1})^\top$ ,  
     $\mathbf{b}_t = \mathbf{b}_{t-1} + \hat{r}_t \boldsymbol{\phi}(\mathbf{x}_{t,a_t}; \mathbf{w}_{t-1})$ ,
  - 7:   update  $\boldsymbol{\theta}_t = \mathbf{A}_t^{-1} \mathbf{b}_t$
  - 8:   **if**  $\operatorname{mod}(t, H) = 0$  **then**
  - 9:      $\mathbf{w}_t \leftarrow$  output of Algorithm 2
  - 10:     $q = q + 1$
  - 11:   **else**
  - 12:      $\mathbf{w}_t = \mathbf{w}_{t-1}$
  - 13:   **end if**
  - 14: **end for**
  - 15: **Output**  $\mathbf{w}_T$
- 

---

**Algorithm 2** Update Weight Parameters with Gradient Descent

---

- 1: **Input:** initial point  $\mathbf{w}_q^{(0)} = \mathbf{w}^{(0)}$ , maximum iteration number  $n$ , step size  $\eta_q$ , and loss function defined in (3.5).
  - 2: **for**  $s = 1, \dots, n$  **do**
  - 3:    $\mathbf{w}_q^{(s)} = \mathbf{w}_q^{(s-1)} - \eta_q \nabla_{\mathbf{w}} \mathcal{L}_q(\mathbf{w}_q^{(s-1)})$ .
  - 4: **end for**
  - 5: **Output**  $\mathbf{w}_q^{(n)}$
- 

217 **Assumption 4.1.** For all  $i \geq 1$  and  $k \in [K]$ , we assume that  $\|\mathbf{x}_{i,k}\|_2 = 1$  and its entries satisfy  
218  $[\mathbf{x}_{i,k}]_j = [\mathbf{x}_{j,k}]_{j+d/2}$ .

219 The assumption that  $\|\mathbf{x}_{i,k}\|_2 = 1$  is not essential and is only imposed for simplicity, which is also  
220 used in Zou and Gu [53], Zhou et al. [52]. Finally, the condition on the entries of  $\mathbf{x}_{i,k}$  is also mild  
221 since otherwise we could always construct  $\mathbf{x}'_{i,k} = [\mathbf{x}_{i,k}^\top, \mathbf{x}_{i,k}^\top]^\top / \sqrt{2}$  to replace it. An implication of  
222 Assumption 4.1 is that the initialization scheme in Algorithm 1 results in  $\boldsymbol{\phi}(\mathbf{x}_{i,k}; \mathbf{w}^{(0)}) = \mathbf{0}$  for all  
223  $i \in [T]$  and  $k \in [K]$ .

224 We assume the following stability condition on the spectral norm of the neural network gradient:

225 **Assumption 4.2.** There is a constant  $\ell_{\text{Lip}} > 0$  such that it holds

$$\left\| \frac{\partial \boldsymbol{\phi}}{\partial \mathbf{w}}(\mathbf{x}; \mathbf{w}_0) - \frac{\partial \boldsymbol{\phi}}{\partial \mathbf{w}}(\mathbf{x}'; \mathbf{w}_0) \right\|_2 \leq \ell_{\text{Lip}} \|\mathbf{x} - \mathbf{x}'\|_2,$$

226 for all  $\mathbf{x}, \mathbf{x}' \in \{\mathbf{x}_{i,k}\}_{i \in [T], k \in [K]}$ .

227 The inequality in Assumption 4.2 resembles the Lipschitz condition on the gradient of the neural  
228 network. However, it is essentially different from the smoothness condition since here the gradient  
229 is taken with respect to the neural network weights while the Lipschitz condition is imposed on the  
230 feature parameter  $\mathbf{x}$ . Similar conditions are widely made in nonconvex optimization [46, 10, 48], in  
231 the name of first-order stability, which is essential to derive the convergence of alternating optimization  
232 algorithms. Furthermore, Assumption 4.2 is only required on the  $TK$  training data points and a  
233 specific weight parameter  $\mathbf{w}_0$ . Therefore, the condition will hold if the raw feature data lie in a  
234 certain subspace of  $\mathbb{R}^d$ . We provided some further discussions in the supplementary material about  
235 this assumption for interested readers.



236 In order to analyze the regret bound of Algorithm 1, we need to characterize the properties of the  
 237 deep neural network in (2.2) that is used to represent the feature vectors. Following a recent line of  
 238 research [27, 12, 7, 52], we define the covariance between two data point  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$  as follows.

$$\begin{aligned}\tilde{\Sigma}^{(0)}(\mathbf{x}, \mathbf{y}) &= \Sigma^{(0)}(\mathbf{x}, \mathbf{y}) = \mathbf{x}^\top \mathbf{y}, \\ \Lambda^{(l)}(\mathbf{x}, \mathbf{y}) &= \begin{bmatrix} \Sigma^{l-1}(\mathbf{x}, \mathbf{x}) & \Sigma^{l-1}(\mathbf{x}, \mathbf{y}) \\ \Sigma^{l-1}(\mathbf{y}, \mathbf{x}) & \Sigma^{l-1}(\mathbf{y}, \mathbf{y}) \end{bmatrix}, \\ \Sigma^{(l)}(\mathbf{x}, \mathbf{y}) &= 2\mathbb{E}_{(u,v) \sim N(\mathbf{0}, \Lambda^{(l-1)}(\mathbf{x}, \mathbf{y}))}[\sigma(u)\sigma(v)], \\ \tilde{\Sigma}^{(l)}(\mathbf{x}, \mathbf{y}) &= 2\tilde{\Sigma}^{(l-1)}(\mathbf{x}, \mathbf{y})\mathbb{E}_{u,v}[\dot{\sigma}(u)\dot{\sigma}(v)] + \Sigma^{(l)}(\mathbf{x}, \mathbf{y}),\end{aligned}\tag{4.1}$$

239 where  $(u, v) \sim N(\mathbf{0}, \Lambda^{(l-1)}(\mathbf{x}, \mathbf{y}))$ , and  $\dot{\sigma}(\cdot)$  is the derivative of activation function  $\sigma(\cdot)$ . We denote  
 240 the neural tangent kernel (NTK) matrix  $\mathbf{H} \in \mathbb{R}^{TK \times TK}$  based on all feature vectors  $\{\mathbf{x}_{t,k}\}_{t \in [T], k \in [K]}$ .  
 241 Renumbering  $\{\mathbf{x}_{t,k}\}_{t \in [T], k \in [K]}$  as  $\{\mathbf{x}_i\}_{i=1, \dots, TK}$ , then each entry  $\mathbf{H}_{ij}$  is defined as

$$\mathbf{H}_{ij} = \frac{1}{2}(\tilde{\Sigma}^{(L)}(\mathbf{x}_i, \mathbf{x}_j) + \Sigma^{(L)}(\mathbf{x}_i, \mathbf{x}_j)),\tag{4.2}$$

242 for all  $i, j \in [TK]$ . Based on the above definition, we impose the following assumption on  $\mathbf{H}$ .

243 **Assumption 4.3.** The neural tangent kernel defined in (4.2) is positive definite, i.e.,  $\lambda_{\min}(\mathbf{H}) \geq \lambda_0$   
 244 for some constant  $\lambda_0 > 0$ .

245 Assumption 4.3 essentially requires the neural tangent kernel matrix  $\mathbf{H}$  to be non-singular, which is  
 246 a mild condition and also imposed in other related work [21, 7, 12, 52]. Moreover, it is shown that  
 247 Assumption 4.3 can be easily derived from Assumption 4.1 for two-layer ReLU networks [37, 53].  
 248 Therefore, Assumption 4.3 is mild or even negligible given the non-degeneration assumption on the  
 249 feature vectors. Also note that matrix  $\mathbf{H}$  is only defined based on layers  $l = 1, \dots, L$  of the neural  
 250 network, and does not depend on the output layer  $\theta$ . It is easy to extend the definition of  $\mathbf{H}$  to the  
 251 NTK matrix defined on all layers including the output layer  $\theta$ , which would also be positive definite  
 252 by Assumption 4.3 and the recursion in (4.2).

253 Before we present the regret analysis of the neural contextual bandit, we need to modify the regret  
 254 defined in (2.1) to account for the randomness of the neural network initialization. For a fixed time  
 255 horizon  $T$ , we define the regret of Algorithm 1 as follows.

$$R_T = \mathbb{E} \left[ \sum_{t=1}^T (\hat{r}(\mathbf{x}_{t, a_t^*}) - \hat{r}(\mathbf{x}_{t, a_t})) | \mathbf{w}^{(0)} \right],\tag{4.3}$$

256 where the expectation is taken over the randomness of the reward noise. Note that  $R_T$  defined in (4.3)  
 257 is still a random variable since the initialization of Algorithm 2 is randomly generated.

258 Now we are going to present the regret bound of the proposed algorithm.

259 **Theorem 4.4.** Suppose Assumptions 4.1, 4.2 and 4.3 hold. Assume that  $\|\theta^*\|_2 \leq M$  for some  
 260 positive constant  $M > 0$ . For any  $\delta \in (0, 1)$ , let us choose  $\alpha_t$  in Neural-LinUCB as

$$\alpha_t = \nu \sqrt{2(d \log(1 + t \log(HK)/\lambda) + \log(1/\delta))} + \lambda^{1/2} M.$$

261 We choose the step size  $\eta_q$  of Algorithm 2 as

$$\eta_q \leq C_0 (d^2 mn T^{5.5} L^6 \log(TK/\delta))^{-1},$$

262 and the width of the neural network satisfies  $m = \text{poly}(L, d, 1/\delta, H, \log(TK/\delta))$ . With probability  
 263 at least  $1 - \delta$  over the randomness of the initialization of the neural network, it holds that

$$R_T \leq C_1 \alpha_T \sqrt{Td \log \left( 1 + \frac{TG^2}{\lambda d} \right)} + \frac{C_2 \ell_{\text{Lip}} L^3 d^{5/2} T \sqrt{\log m \log(\frac{1}{\delta}) \log(\frac{TK}{\delta})} \|\mathbf{r} - \tilde{\mathbf{r}}\|_{\mathbf{H}^{-1}}}{m^{1/6}},$$

264 where  $\{C_i\}_{i=0,1,2}$  are absolute constants independent of the problem parameters,  $\mathbf{r} =$   
 265  $(r(\mathbf{x}_1), r(\mathbf{x}_2), \dots, r(\mathbf{x}_{TK}))^\top \in \mathbb{R}^{TK}$  and  $\tilde{\mathbf{r}} = (f(\mathbf{x}_1; \theta_0, \mathbf{w}_0), \dots, f(\mathbf{x}_{TK}; \theta_{T-1}, \mathbf{w}_{T-1}))^\top \in$   
 266  $\mathbb{R}^{TK}$ , and  $\|\mathbf{r}\|_{\mathbf{A}} = \sqrt{\mathbf{r}^\top \mathbf{A} \mathbf{r}}$ .

267 **Remark 4.5.** Theorem 4.4 shows that the regret of Algorithm 1 can be bounded by two parts: the  
 268 first part is of order  $\tilde{O}(\sqrt{T})$ , which resembles the regret bound of linear contextual bandits [1]; the  
 269 second part is of order  $\tilde{O}(m^{-1/6}T\sqrt{(\mathbf{r} - \tilde{\mathbf{r}})^\top \mathbf{H}^{-1}(\mathbf{r} - \tilde{\mathbf{r}})})$ , which depends on the estimation error  
 270 of the neural network  $f$  for the reward generating function  $r$  and the neural tangent kernel  $\mathbf{H}$ .

271 It is worth noting that our theoretical analysis depends on the reward structure assumption that  
 272  $r(\cdot) = \langle \boldsymbol{\theta}_*, \boldsymbol{\psi}(\cdot) \rangle$ . However, the linear structure between  $\boldsymbol{\theta}_*$  and  $\boldsymbol{\psi}(\cdot)$  is not essential. As long as  
 273 the deep representation of the feature vector and the uncertainty weight parameter can be decoupled,  
 274 Algorithm 1 can be easily extended to settings with milder assumptions on the reward structure  
 275 such as generalized linear models [41, 24, 35, 28]. For more general bandit models where no  
 276 assumption is imposed to the reward generating function, it is still unclear whether the decoupled  
 277 deep representation and shallow exploration would work especially in cases a thorough exploration  
 278 may be needed.

279 Based on the result in Theorem 4.4, we can easily verify the following conclusion:

280 **Corollary 4.6.** Under the same conditions of Theorem 4.4, if we choose a sufficiently overpa-  
 281 rameterized neural network mapping  $\phi(\cdot)$  such that  $m \geq T^3$ , then the regret of Algorithm 1 is  
 282  $R_T = \tilde{O}(\sqrt{T}\sqrt{(\mathbf{r} - \tilde{\mathbf{r}})^\top \mathbf{H}^{-1}(\mathbf{r} - \tilde{\mathbf{r}})})$ .

283 **Remark 4.7.** For the ease of presentation, let us denote  $\mathcal{E} := \|\mathbf{r} - \tilde{\mathbf{r}}\|_{\mathbf{H}^{-1}}$ . If we have  $\mathcal{E} = O(1)$ , the  
 284 total regret in Theorem 4.4 becomes  $\tilde{O}(\sqrt{T})$  which matches the regret of linear contextual bandits  
 285 [1]. We remark that there is a similar assumption in [52] where they assume that  $\mathbf{r}^\top \mathbf{H}^{-1} \mathbf{r}$  can be  
 286 upper bounded by a constant. They show that this term can be bounded by the RKHS norm of  $\mathbf{r}$  if  
 287 it belongs to the RKHS induced by the neural tangent kernel [6, 7, 33]. In addition,  $\mathcal{E}$  here is the  
 288 difference between the true reward function and the neural network function, which can also be small  
 289 if the deep neural network function well approximates the reward generating function  $r(\cdot)$ .

## 290 5 Experiments

291 In this section, we provide empirical evaluations of Neural-LinUCB on real-world datasets. As  
 292 we have discussed in Section 3, Neural-LinUCB could be viewed as an instantiation of the Neural-  
 293 Linear scheme studied in Riquelme et al. [38] except that we use the UCB exploration instead of the  
 294 posterior sampling exploration therein. Note that there has been an extensive comparison [38] of the  
 295 Neural-Linear methods with many other baselines such as greedy algorithms, Variational Inference,  
 296 Expectation-Propagation, Bayesian Non-parametrics and so on. Therefore, we do not seek a thorough  
 297 empirical comparison of Neural-LinUCB with all existing bandits algorithms. We refer readers who  
 298 are interested in the performance of Neural-Linear methods with deep representation and shallow  
 299 exploration compared with a vast of baselines in the literature to the benchmark study by Riquelme  
 300 et al. [38]. In this experiment, we only aim to show the advantages of our algorithm over the following  
 301 baselines: (1) Neural-Linear [38]; (2) LinUCB [16], which does not have a *deep representation* of the  
 302 feature vectors; and (3) NeuralUCB [52], which performs UCB exploration on all the parameters of  
 303 the neural network instead of the *shallow exploration* used in our paper. All numerical experiments  
 304 were run on a workstation with Intel(R) Xeon(R) CPU E5-2637 v4 @ 3.50GHz.

305 **Datasets:** we evaluate the performances of all algorithms on bandit problems created from real-world  
 306 data. Specifically, following the experimental setting in Zhou et al. [52], we use datasets (*Shuttle*)  
 307 *Statlog*, *Magic* and *Coverttype* from UCI machine learning repository [23], and the *MINST* dataset  
 308 from LeCun et al. [31]. The details of these datasets are presented in Table 1. In Table 1, each  
 309 instance represents a feature vector  $\mathbf{x} \in \mathbb{R}^d$  that is associated with one of the  $K$  arms, and dimension  
 310  $d$  is the number of attributes in each instance.

Table 1: Specifications of datasets from the UCI machine learning repository used in this paper.

	<i>Statlog</i>	<i>Magic</i>	<i>Coverttype</i>	<i>MNIST</i>
Number of attributes	9	11	54	784
Number of arms	7	2	7	10
Number of instances	58,000	19,020	581,012	60,000



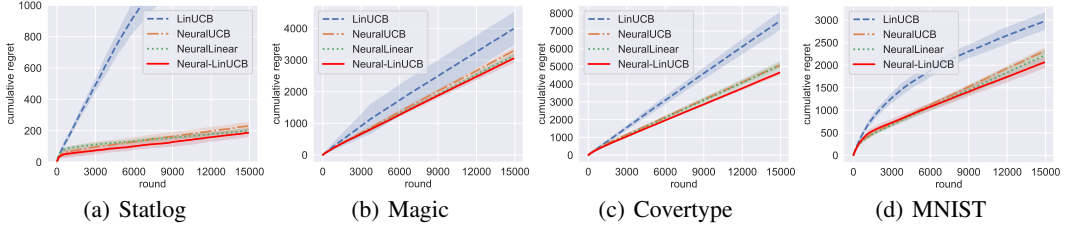


Figure 1: The cumulative regrets of LinUCB, NeuralUCB, Neural-Linear and Neural-LinUCB over 15,000 rounds. Experiments are averaged over 10 repetitions.

311 **Implementations:** for LinUCB, we follow the setting in Li et al. [34] to use disjoint models  
 312 for different arms. For neural network based algorithms such as NeuralUCB, Neural-Linear and  
 313 Neural-LinUCB, we use a ReLU neural network defined as in (2.2) with  $L = 2$  and 2000 for the  
 314 UCI datasets (*Statlog*, *Magic*, *Coverttype*). Thus the neural network weights are  $\mathbf{W}_1 \in \mathbb{R}^{m \times d}$ ,  
 315  $\mathbf{W}_2 \in \mathbb{R}^{k \times m}$ , and  $\boldsymbol{\theta} \in \mathbb{R}^k$  respectively, where  $k = 100$ ,  $m = 2000$ , and  $d$  is the dimension of  
 316 features in the corresponding task. Since the problem size of the MNIST dataset is larger, inspired  
 317 by Hinton and Salakhutdinov [26], we use a deeper NN and set  $L = 3$ ,  $k = 100$  and  $m = 100$ ,  
 318 with weights  $\mathbf{W}_1 \in \mathbb{R}^{m \times d}$ ,  $\mathbf{W}_2 \in \mathbb{R}^{m \times m}$ ,  $\mathbf{W}_3 \in \mathbb{R}^{k \times m}$ , and  $\boldsymbol{\theta} \in \mathbb{R}^k$ . We set the time horizon  
 319  $T = 15,000$ , which is the total number of rounds for each algorithm on each dataset. We use  
 320 gradient decent to optimize the network weights, with a step size  $\eta_g = 1e-5$  and maximum iteration  
 321 number  $n = 1,000$ . To speed up the training process, the network parameter  $\mathbf{w}$  is updated every  
 322  $H = 100$  rounds starting from round 2000. We also apply early stopping when the loss difference  
 323 of two consecutive iterations is smaller than a threshold of  $1e-6$ . We set  $\lambda = 1$  and  $\alpha_t = 0.02$   
 324 for all algorithms,  $t \in [T]$ . Following the setting in Riquelme et al. [38], we use round-robin to  
 325 independently select each arm for 3 times at the beginning of each algorithm. For NeuralUCB, since  
 326 it is computationally unaffordable to perform the original UCB exploration as displayed in Zhou et al.  
 327 [52], we follow their experimental setting to replace the matrix  $\mathbf{Z}_t \in \mathbb{R}^{(d+\bar{p}) \times (d+\bar{p})}$  in Zhou et al.  
 328 [52] with its diagonal matrix.

329 **Results:** we plot the cumulative regret of all algorithms versus round in Figures 1(a), 1(b) and 1(c)  
 330 for UCI datasets and in Figure 1(d) for MNIST. The results are reported based on the average of  
 331 10 repetitions over different random shuffles of the datasets. It can be seen that algorithms based  
 332 on neural network representations (NeuralUCB, Neural-Linear and Neural-LinUCB) consistently  
 333 outperform the linear contextual bandit method LinUCB, which shows that linear models may  
 334 lack representation power and find biased estimates for the underlying reward generating function.  
 335 Furthermore, our proposed Neural-LinUCB achieves a comparable regret with NeuralUCB in all  
 336 experiments despite the fact that our algorithm only explores in the output layer of the neural network,  
 337 which is more computationally efficient as we will show in the sequel. The results in our experiment  
 338 are well aligned with our theory that deep representation and shallow exploration are sufficient to  
 339 guarantee a good performance of neural contextual bandit algorithms, which is also consistent with  
 340 the findings in existing literature [38] that decoupling the representation learning and uncertainty  
 341 estimation improves the performance.

342 We also conducted experiments to study the effects of different widths of deep neural networks on  
 343 the regret performance and to show the computational efficiency of Neural-LinUCB compared with  
 344 existing neural bandit algorithms. Due to the space limit, we defer the results to Appendix A.

## 345 6 Conclusions

346 In this paper, we propose a new neural contextual bandit algorithm called Neural-LinUCB, which uses  
 347 the hidden layers of a ReLU neural network as a deep representation of the raw feature vectors and  
 348 performs UCB type exploration on the last layer of the neural network. By incorporating techniques  
 349 in liner contextual bandits and neural tangent kernels, we prove that the proposed algorithm achieves  
 350 a sublinear regret when the width of the network is sufficiently large. This is the first regret analysis  
 351 of neural contextual bandit algorithms with deep representation and shallow exploration, which have  
 352 been observed in practice to work well on many benchmark bandit problems [38]. We also conducted  
 353 experiments on real-world datasets to demonstrate the advantage of the proposed algorithm over  
 354 LinUCB and existing neural contextual bandit algorithms.

## References

- [1] Yasin Abbasi-Yadkori, Dávid Pál, and Csaba Szepesvári. Improved algorithms for linear stochastic bandits. In *Advances in Neural Information Processing Systems*, pages 2312–2320, 2011.
- [2] Alekh Agarwal, Daniel Hsu, Satyen Kale, John Langford, Lihong Li, and Robert Schapire. Taming the monster: A fast and simple algorithm for contextual bandits. In *International Conference on Machine Learning*, pages 1638–1646, 2014.
- [3] Zeyuan Allen-Zhu, Yuanzhi Li, and Yingyu Liang. Learning and generalization in overparameterized neural networks, going beyond two layers. In *Advances in neural information processing systems*, pages 6155–6166, 2019.
- [4] Zeyuan Allen-Zhu, Yuanzhi Li, and Zhao Song. A convergence theory for deep learning via over-parameterization. In *International Conference on Machine Learning*, pages 242–252, 2019.
- [5] Robin Allesiardo, Raphaël Féraud, and Djallel Bouneffouf. A neural networks committee for the contextual bandit problem. In *International Conference on Neural Information Processing*, pages 374–381. Springer, 2014.
- [6] Sanjeev Arora, Simon Du, Wei Hu, Zhiyuan Li, and Ruosong Wang. Fine-grained analysis of optimization and generalization for overparameterized two-layer neural networks. In *International Conference on Machine Learning*, pages 322–332, 2019.
- [7] Sanjeev Arora, Simon S Du, Wei Hu, Zhiyuan Li, Russ R Salakhutdinov, and Ruosong Wang. On exact computation with an infinitely wide neural net. In *Advances in Neural Information Processing Systems*, pages 8139–8148, 2019.
- [8] Jean-Yves Audibert, Rémi Munos, and Csaba Szepesvári. Exploration–exploitation tradeoff using variance estimates in multi-armed bandits. *Theoretical Computer Science*, 410(19): 1876–1902, 2009.
- [9] Peter Auer, Nicolo Cesa-Bianchi, and Paul Fischer. Finite-time analysis of the multiarmed bandit problem. *Machine learning*, 47(2-3):235–256, 2002.
- [10] Sivaraman Balakrishnan, Martin J Wainwright, Bin Yu, et al. Statistical guarantees for the em algorithm: From population to sample-based analysis. *The Annals of Statistics*, 45(1):77–120, 2017.
- [11] Qi Cai, Zhuoran Yang, Jason D Lee, and Zhaoran Wang. Neural temporal-difference learning converges to global optima. In *Advances in Neural Information Processing Systems*, 2019.
- [12] Yuan Cao and Quanquan Gu. A generalization theory of gradient descent for learning over-parameterized deep relu networks. *arXiv preprint arXiv:1902.01384*, 2019.
- [13] Yuan Cao and Quanquan Gu. Generalization bounds of stochastic gradient descent for wide and deep neural networks. In *Advances in Neural Information Processing Systems*, pages 10835–10845, 2019.
- [14] Olivier Chapelle and Lihong Li. An empirical evaluation of thompson sampling. In *Advances in neural information processing systems*, pages 2249–2257, 2011.
- [15] Sayak Ray Chowdhury and Aditya Gopalan. On kernelized multi-armed bandits. In *International Conference on Machine Learning*, pages 844–853, 2017.
- [16] Wei Chu, Lihong Li, Lev Reyzin, and Robert Schapire. Contextual bandits with linear payoff functions. In *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*, pages 208–214, 2011.
- [17] Mark Collier and Hector Urdiales Llorens. Deep contextual multi-armed bandits. *arXiv preprint arXiv:1807.09809*, 2018.

- 401 [18] George Cybenko. Approximation by superpositions of a sigmoidal function. *Mathematics of*  
402 *control, signals and systems*, 2(4):303–314, 1989.
- 403 [19] Varsha Dani, Thomas P Hayes, and Sham M Kakade. Stochastic linear optimization under  
404 bandit feedback. In *Conference on Learning Theory*, 2008.
- 405 [20] Aniket Anand Deshmukh, Abhimanu Kumar, Levi Boyles, Denis Charles, Eren Manavoglu,  
406 and Urun Dogan. Self-supervised contextual bandits in computer vision. *arXiv preprint*  
407 *arXiv:2003.08485*, 2020.
- 408 [21] Simon Du, Jason Lee, Haochuan Li, Liwei Wang, and Xiyu Zhai. Gradient descent finds global  
409 minima of deep neural networks. In *International Conference on Machine Learning*, pages  
410 1675–1685, 2019.
- 411 [22] Simon S. Du, Xiyu Zhai, Barnabas Poczos, and Aarti Singh. Gradient descent provably  
412 optimizes over-parameterized neural networks. In *International Conference on Learning*  
413 *Representations*, 2019. URL <https://openreview.net/forum?id=SlE3i09YQ>.
- 414 [23] Dheeru Dua and Casey Graff. UCI machine learning repository, 2017. URL [http://archive.](http://archive.ics.uci.edu/ml)  
415 [ics.uci.edu/ml](http://archive.ics.uci.edu/ml).
- 416 [24] Sarah Filippi, Olivier Cappé, Aurélien Garivier, and Csaba Szepesvári. Parametric bandits: The  
417 generalized linear case. In *Advances in Neural Information Processing Systems*, pages 586–594,  
418 2010.
- 419 [25] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep learning*. MIT press, 2016.
- 420 [26] Geoffrey E Hinton and Ruslan R Salakhutdinov. Reducing the dimensionality of data with  
421 neural networks. *science*, 313(5786):504–507, 2006.
- 422 [27] Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and  
423 generalization in neural networks. In *Advances in neural information processing systems*, pages  
424 8571–8580, 2018.
- 425 [28] Branislav Kveton, Manzil Zaheer, Csaba Szepesvari, Lihong Li, Mohammad Ghavamzadeh,  
426 and Craig Boutilier. Randomized exploration in generalized linear bandits. In *International*  
427 *Conference on Artificial Intelligence and Statistics*, pages 2066–2076, 2020.
- 428 [29] John Langford and Tong Zhang. The epoch-greedy algorithm for multi-armed bandits with side  
429 information. In *Advances in neural information processing systems*, pages 817–824, 2008.
- 430 [30] Tor Lattimore and Csaba Szepesvári. *Bandit algorithms*. Cambridge University Press, 2020.
- 431 [31] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning  
432 applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- 433 [32] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *nature*, 521(7553):436–444,  
434 2015.
- 435 [33] Jaehoon Lee, Lechao Xiao, Samuel Schoenholz, Yasaman Bahri, Roman Novak, Jascha Sohl-  
436 Dickstein, and Jeffrey Pennington. Wide neural networks of any depth evolve as linear models  
437 under gradient descent. In *Advances in neural information processing systems*, pages 8570–8581,  
438 2019.
- 439 [34] Lihong Li, Wei Chu, John Langford, and Robert E Schapire. A contextual-bandit approach to  
440 personalized news article recommendation. In *Proceedings of the 19th international conference*  
441 *on World wide web*, pages 661–670, 2010.
- 442 [35] Lihong Li, Yu Lu, and Dengyong Zhou. Provably optimal algorithms for generalized linear  
443 contextual bandits. In *International Conference on Machine Learning*, pages 2071–2080, 2017.
- 444 [36] Boyi Liu, Qi Cai, Zhuoran Yang, and Zhaoran Wang. Neural trust region/proximal policy  
445 optimization attains globally optimal policy. In *Advances in Neural Information Processing*  
446 *Systems*, pages 10564–10575, 2019.

- 447 [37] Samet Oymak and Mahdi Soltanolkotabi. Towards moderate overparameterization: global  
448 convergence guarantees for training shallow neural networks. *IEEE Journal on Selected Areas*  
449 *in Information Theory*, 2020.
- 450 [38] Carlos Riquelme, George Tucker, and Jasper Snoek. Deep bayesian bandits showdown: An  
451 empirical comparison of bayesian deep networks for thompson sampling. In *International*  
452 *Conference on Learning Representations*, 2018. URL [https://openreview.net/forum?](https://openreview.net/forum?id=SyYe6k-CW)  
453 [id=SyYe6k-CW](https://openreview.net/forum?id=SyYe6k-CW).
- 454 [39] Paat Rusmevichientong and John N Tsitsiklis. Linearly parameterized bandits. *Mathematics of*  
455 *Operations Research*, 35(2):395–411, 2010.
- 456 [40] Daniel J. Russo, Benjamin Van Roy, Abbas Kazerouni, Ian Osband, and Zheng Wen. A tutorial  
457 on thompson sampling. *Foundations and Trends® in Machine Learning*, 11(1):1–96, 2018.  
458 ISSN 1935-8237.
- 459 [41] Jyotirmoy Sarkar. One-armed bandit problems with covariates. *The Annals of Statistics*, pages  
460 1978–2002, 1991.
- 461 [42] Jasper Snoek, Oren Rippel, Kevin Swersky, Ryan Kiros, Nadathur Satish, Narayanan Sundaram,  
462 Mostofa Patwary, Mr Prabhat, and Ryan Adams. Scalable bayesian optimization using deep  
463 neural networks. In *International conference on machine learning*, pages 2171–2180, 2015.
- 464 [43] William R Thompson. On the likelihood that one unknown probability exceeds another in view  
465 of the evidence of two samples. *Biometrika*, 25(3/4):285–294, 1933.
- 466 [44] Michal Valko, Nathan Korda, Rémi Munos, Ilias Flaounas, and Nello Cristianini. Finite-time  
467 analysis of kernelised contextual bandits. In *Proceedings of the Twenty-Ninth Conference on*  
468 *Uncertainty in Artificial Intelligence*, pages 654–663, 2013.
- 469 [45] Lingxiao Wang, Qi Cai, Zhuoran Yang, and Zhaoran Wang. Neural policy gradient meth-  
470 ods: Global optimality and rates of convergence. In *International Conference on Learning*  
471 *Representations*, 2020. URL <https://openreview.net/forum?id=BJgQfkSYDS>.
- 472 [46] Zhaoran Wang, Han Liu, and Tong Zhang. Optimal computational and statistical rates of  
473 convergence for sparse nonconvex learning problems. *Annals of statistics*, 42(6):2164, 2014.
- 474 [47] Pan Xu and Quanquan Gu. A finite-time analysis of q-learning with neural network function  
475 approximation. In *International Conference on Machine Learning*, 2020.
- 476 [48] Pan Xu, Jian Ma, and Quanquan Gu. Speeding up latent variable gaussian graphical model  
477 estimation via nonconvex optimization. In *Advances in Neural Information Processing Systems*,  
478 pages 1933–1944, 2017.
- 479 [49] Tom Zahavy and Shie Mannor. Deep neural linear bandits: Overcoming catastrophic forgetting  
480 through likelihood matching. *arXiv preprint arXiv:1901.08612*, 2019.
- 481 [50] Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding  
482 deep learning requires rethinking generalization. In *International Conference on Learning*  
483 *Representations*, 2017. URL <https://openreview.net/forum?id=Sy8gdB9xx>.
- 484 [51] Weitong Zhang, Dongruo Zhou, Lihong Li, and Quanquan Gu. Neural thompson sampling.  
485 *arXiv preprint arXiv:2010.00827*, 2020.
- 486 [52] Dongruo Zhou, Lihong Li, and Quanquan Gu. Neural contextual bandits with ucb-based  
487 exploration. In *International Conference on Machine Learning*, 2020.
- 488 [53] Difan Zou and Quanquan Gu. An improved analysis of training over-parameterized deep neural  
489 networks. In *Advances in Neural Information Processing Systems*, pages 2053–2062, 2019.
- 490 [54] Difan Zou, Yuan Cao, Dongruo Zhou, and Quanquan Gu. Stochastic gradient descent optimizes  
491 over-parameterized deep relu networks. *arXiv preprint arXiv:1811.08888*, 2018.

492 **Checklist**

- 493 1. For all authors...
- 494 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s  
495 contributions and scope? [Yes]
- 496 (b) Did you describe the limitations of your work? [Yes] We discussed the limitation of  
497 the assumptions made in this paper. We also admit in the experiment that the theory  
498 maybe conservative since our experiment does not require a very wide neural network  
499 to achieve good performance.
- 500 (c) Did you discuss any potential negative societal impacts of your work? [N/A] This work  
501 focuses on a general methodology in bandit problems and its theoretical analysis. It  
502 does not cause any negative social impact.
- 503 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
504 them? [Yes]
- 505 2. If you are including theoretical results...
- 506 (a) Did you state the full set of assumptions of all theoretical results? [Yes] See the  
507 assumptions listed in Section 4
- 508 (b) Did you include complete proofs of all theoretical results? [Yes] Proofs are provided in  
509 the appendix.
- 510 3. If you ran experiments...
- 511 (a) Did you include the code, data, and instructions needed to reproduce the main experi-  
512 mental results (either in the supplemental material or as a URL)? [Yes] We provide  
513 them in the supplementary material.
- 514 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
515 were chosen)? [Yes] We specify all the details in the Implementations paragraph of  
516 Section 5.
- 517 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
518 ments multiple times)? [Yes] All the figures are plotted with the standard error with  
519 respect to random repetitions.
- 520 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
521 of GPUs, internal cluster, or cloud provider)? [Yes] We stated the type of workstation  
522 at the end of the first paragraph of Section 5.
- 523 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 524 (a) If your work uses existing assets, did you cite the creators? [Yes] As we mentioned in  
525 Section 5, we used codes from baseline algorithms and public available datasets. All  
526 the assets were properly cited.
- 527 (b) Did you mention the license of the assets? [N/A] All the codes and datasets are  
528 open-source.
- 529 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]  
530 We include our code in the supplementary for reproduction.
- 531 (d) Did you discuss whether and how consent was obtained from people whose data you’re  
532 using/curating? [N/A]
- 533 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
534 information or offensive content? [N/A] The data does not contain any personally  
535 identifiable information or offensive content.
- 536 5. If you used crowdsourcing or conducted research with human subjects...
- 537 (a) Did you include the full text of instructions given to participants and screenshots, if  
538 applicable? [N/A]
- 539 (b) Did you describe any potential participant risks, with links to Institutional Review  
540 Board (IRB) approvals, if applicable? [N/A]
- 541 (c) Did you include the estimated hourly wage paid to participants and the total amount  
542 spent on participant compensation? [N/A]