



Overview

- Contextual Bandits are rich models for online decision-making problems where an agent sequentially interacts with an unknown, responsive environment and receives a reward.
- How can an agent leverage the expressive power of neural networks to learn the reward function and eventually converge to an optimal action selection policy?
- We propose algorithms that employ (convolutional) neural networks to estimate the reward, and provably attain sublinear regret.

Problem Setting

- At every step $1 \leq t \leq T$

$$y_t = f(\mathbf{x}_t) + \epsilon_t$$

Action	$\mathbf{a}_t \in \mathcal{A}$ one-hot vector of length $ \mathcal{A} $
Context	$\mathbf{z}_t = (\mathbf{z}_{t,1}, \dots, \mathbf{z}_{t, \mathcal{A} }) \in \mathbb{R}^{d \times \mathcal{A} }$
Noise	$\epsilon_t : \sigma^2$ sub-Gaussian, i.i.d.
Input	$\mathbf{x}_t = \mathbf{z}_t \mathbf{a}_t \in \mathbb{R}^d \quad \mathcal{X} \subseteq \mathbb{S}^{d-1}$
Reward	$f : \mathcal{X} \rightarrow \mathbb{R} \quad f \in \mathcal{H}_{k_{\text{NN}}}$ $\ f\ _{k_{\text{NN}}} \leq B \quad k_{\text{NN}}$: the Neural Tangent Kernel
Cumulative regret	$R_T = \sum_{t=1}^T f(\mathbf{x}_t^*) - f(\mathbf{x}_t) \quad \mathbf{x}_t^* = \arg \max_{x=z_t a, a \in \mathcal{A}} f(x)$
Learner's goal	$R_T / T \rightarrow 0$ as $T \rightarrow \infty$

Neural Bandits

- We propose two algorithms, NN-UCB and CNN-UCB.
- At every step $1 \leq t \leq T$

Approximate the mean of reward with $\hat{\mu}_{t-1}(\mathbf{x}) = f^{(J)}(\mathbf{x}; \theta_{t-1})$
the (C)NN trained for J steps of gradient descent on

$$\mathcal{L}(\theta) = \sum_{i=1}^{t-1} (f(\mathbf{x}_i; \theta) - y_i)^2 + m\sigma^2 \|\theta - \theta^0\|_2^2$$

Approximate the variance of reward with

$$\hat{\sigma}_{t-1}^2(\mathbf{x}) = \frac{\mathbf{g}^T(\mathbf{x})}{\sqrt{m}} \left(\sigma^2 \mathbf{I} + \sum_{i=1}^{t-1} \frac{\mathbf{g}^T(\mathbf{x}_i) \mathbf{g}(\mathbf{x}_i)}{m} \right)^{-1} \frac{\mathbf{g}(\mathbf{x})}{\sqrt{m}}$$

$$\mathbf{g}(\mathbf{x}) = \nabla_{\theta} f(\mathbf{x}; \theta^{(0)})$$

Pick actions by maximizing the approximate Upper Confidence Bound

$$\mathbf{x}_t = \arg \max_{\mathbf{x}=\mathbf{z}_t \mathbf{a}, \mathbf{a} \in \mathcal{A}} \hat{\mu}_{t-1}(\mathbf{x}) + \sqrt{\beta_t} \hat{\sigma}_{t-1}(\mathbf{x})$$

Information Gain

- It often appears in bandit regret bounds and quantifies the speed at which the agent learns about the reward function. It is the mutual information between the observed reward and true reward values. H denotes the entropy.

$$I(\mathbf{y}_T; \mathbf{f}_T) := H(\mathbf{y}_T) - H(\mathbf{y}_T | \mathbf{f}_T) = \frac{1}{2} \log \det(\mathbf{I} + \sigma^{-2} \mathbf{K}_T)$$

$$\mathbf{K}_T = [k_{\text{NN}}(\mathbf{x}_i, \mathbf{x}_j)]_{i,j \leq T}$$

- Its maximum depends only on the domain, the noise and the kernel function

$$\gamma_T = \max_{\mathbf{x}_1, \dots, \mathbf{x}_T} I(\mathbf{y}_T; \mathbf{f}_T)$$

Theorem (Information Gain Bound, Informal)

Then the maximum information gain associated with the NTK of a fully-connected ReLU network is bounded by

$$\gamma_{T, \text{NN}} = \tilde{\mathcal{O}} \left(C_1(d, L) T^{\frac{d-1}{d}} \right)$$

Main Result

- We resolve the open problem of proving sublinear regret bounds for general context sequences, considering both fully-connected and convolutional nets.

We show that our algorithms converge to the optimal policy in polynomial time with high probability.

Theorem (Regret Bound, Neural, Informal)

Let $\delta \in (1, 0)$. Set $\beta_t = 2 \log(2T|\mathcal{A}|/\delta)$. Choose the width s.t.

$$m \geq \text{poly}(T, L, |\mathcal{A}|, \sigma^{-2}, B, \lambda_0^{-1}, \log(1/\delta))$$

minimum eigenvalue of the kernel matrix

and learning rate $\eta = C(LmT + m\sigma^2)^{-1}$ with some constant C . Then with probability greater than $1 - \delta$, SupNN-UCB satisfies

$$R(T) = \tilde{\mathcal{O}}(C_2(d, L) T^{\frac{2d-1}{2d}}).$$

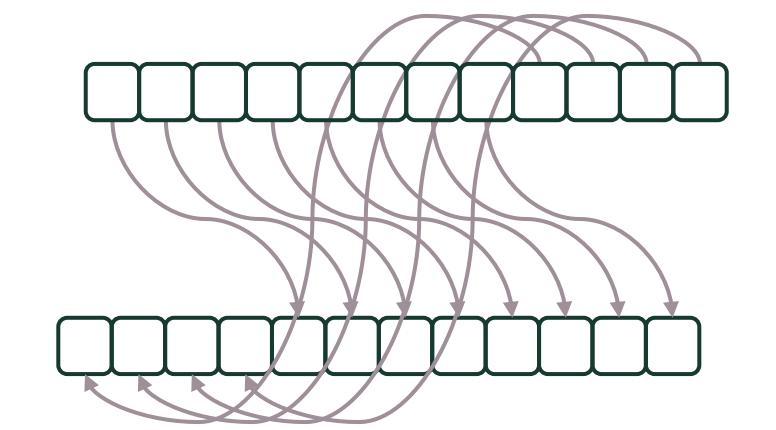
- Comparison to prior work: [2] are the first to introduce neural contextual bandits. They provide a guarantee roughly of the form

$$R_T \leq \tilde{\mathcal{O}}(I(\mathbf{y}_T; \mathbf{f}_T) \sqrt{T})$$

This bound is data-dependent via the Information Gain. The maximum information gain itself grows as $\gamma_{T, \text{NN}} = \tilde{\mathcal{O}}(T^{\frac{d-1}{d}})$ for the NTK. Without further assumptions on the sequence of contexts, this bound may be vacuous.

Convolutional Guarantees

- Invariance Trick [1] $\mathbf{w} * \mathbf{x} = \sum_{l=1}^d \langle \mathbf{w}, c_l \cdot \mathbf{x} \rangle$
 $c_l \cdot \mathbf{x} = (x_{l+1}, x_{l+2}, \dots, x_d, x_1, \dots, x_l)$



2-Layer convolutional network is invariant to circular shifts.

$$f_{\text{CNN}}(\mathbf{x}; \boldsymbol{\theta}) = \frac{1}{d} \sum_{i=1}^m v_i \sigma_{\text{relu}}(\mathbf{w}_i * \mathbf{x}) = \frac{1}{d} \sum_{l=1}^d f_{\text{NN}}(c_l \cdot \mathbf{x}; \boldsymbol{\theta})$$

$$f_{\text{CNN}}(c_l \cdot \mathbf{x}) = f_{\text{CNN}}(\mathbf{x})$$

- Connection between the NTK and the CNTK

$$k_{\text{CNN}}(\mathbf{x}, \mathbf{x}') = \frac{1}{d} \sum_{l=1}^d k_{\text{NN}}(\mathbf{x}, c_l \cdot \mathbf{x}')$$

When restricted to S^{d-1} , the NTK is described by $(\mu_k, \mathcal{F}_{d,k})_{k \geq 0}$ a sequence of eigenvalue eigenspace pairs. Here, the k -th eigenspace is spanned by degree- k polynomials. For the CNTK however, the k -th eigenspace is spanned by polynomials that are *invariant to circular shifts*. Searching through this shrunken RKHS gives us improved rates for CNN-UCB.

- Reward model assumption

$$f \in \mathcal{H}_{k_{\text{CNN}}} \quad \|f\|_{k_{\text{CNN}}} \leq B$$

CNN-UCB benefits from a shift-invariant structure.

Our results suggest that for a high-dimensional input it outperforms NN-UCB.

- Guarantees for the Convolutional Neural Bandit

$$\gamma_{T, \text{CNN}} = \tilde{\mathcal{O}} \left(C_1(d, L) \left(\frac{T}{d} \right)^{\frac{d-1}{d}} \right)$$

Theorem (Regret Bound, Convolutional, Informal)

Let $\delta \in (1, 0)$. Set $\beta_t = 2 \log(2T|\mathcal{A}|/\delta)$. For any T , there exists width m such that if $\eta = C(LmT + m\sigma^2)^{-1}$ with some constant C , then with probability greater than $1 - \delta$, SupCNN-UCB satisfies

$$R(T) = \tilde{\mathcal{O}} \left(\frac{C_2(d, L)}{d^{(d-1)/2d}} T^{\frac{2d-1}{2d}} \right).$$

References

- [1] Song Mei, Theodor Misiakiewicz, and Andrea Montanari. Learning with invariances in random features and kernel models. In Mikhail Belkin and Samory Kpotufe, editors, *Proceedings of Thirty Fourth Conference on Learning Theory*, volume 134 of *Proceedings of Machine Learning Research*, pages 3351–3418. PMLR, 15–19 Aug 2021.
- [2] Dongruo Zhou, Lihong Li, and Quanquan Gu. Neural contextual bandits with ucb-based exploration. In *International Conference on Machine Learning*, pages 11492–11502. PMLR, 2020.