

Neural Contextual Bandits with UCB-based Exploration

Dongruo Zhou ¹ Lihong Li ² Quanquan Gu ¹



¹Department of Computer Science, UCLA

²Google Research

Outline

- | Background
 - | Contextual bandit problem
 - | Deep neural networks

Outline

- | Background
 - | Contextual bandit problem
 - | Deep neural networks
- | Algorithm – NeuralUCB
 - | Use a neural network to learn the reward
 - | Use neural network's gradient to explore
 - | Upper confidence bound strategy

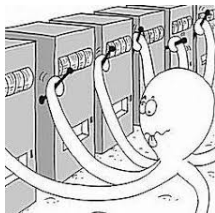
Outline

- | Background
 - | Contextual bandit problem
 - | Deep neural networks
- | Algorithm – NeuralUCB
 - | Use a neural network to learn the reward
 - | Use neural network's gradient to explore
 - | Upper confidence bound strategy
- | Main theory
 - | Neural tangent kernel matrix and effective dimension
 - | $\mathcal{O}(\sqrt{T})$ regret

Background – decision-making problems

Decision-making problems are everywhere!

- | As a gambler in a casino, find a slot machine, you will...
 - | Limited budget, maximize the payoff !
 - | Which arm to pull?
- | As a movie recommender, you need to...
 - | Recommend movies based on users' interests, maximize users' purchase rate
 - | Which movie to recommend?



(a) Slot machine



(b) Movie recommendation

Background – contextual bandit

K -armed contextual bandit problem: movie recommendation

Background – contextual bandit

K -armed contextual bandit problem: movie recommendation

At round t ,

- | Agent observes K d -dimensional contextual vectors (user's movie purchase history)

$$f_{\mathbf{x}_{t,a}} \in \mathbb{R}^d \quad j, a \in [K]$$

Background – contextual bandit

K -armed contextual bandit problem: movie recommendation

At round t ,

- | Agent observes K d -dimensional contextual vectors (user's movie purchase history)

$$f_{\mathbf{x}_{t;a}} \in \mathbb{R}^d \quad j \in [K]$$

- | Agent selects an action a_t and receives a reward $r_{t;a_t}$ (recommends some movie and user choose to purchase or not)

Background – contextual bandit

K -armed contextual bandit problem: movie recommendation

At round t ,

- | Agent observes K d -dimensional contextual vectors (user's movie purchase history)

$$f: \mathbf{x}_{t;a} \in \mathbb{R}^d \rightarrow a \in [K]$$

- | Agent selects an action a_t and receives a reward $r_{t;a_t}$ (recommends some movie and user choose to purchase or not)
- | The goal is to minimize the following pseudo regret

$$R_T = \mathbb{E} \sum_{t=1}^T (r_{t;a_t^*} - r_{t;a_t})$$

where $a_t = \operatorname{argmax}_{a \in [K]} \mathbb{E}[r_{t;a}]$ is the optimal action at round t

Background – contextual linear bandit

$$r_{t,a_t} = h(\mathbf{x}_{t,a_t}) + \epsilon_t, \quad \epsilon_t \text{ -sub-Gaussian}$$

Background – contextual linear bandit

$$r_{t;a_t} = h(\mathbf{x}_{t;a_t})^\top \theta + \epsilon_t, \quad \epsilon_t \text{ -sub-Gaussian}$$

- | Build confidence set for θ and use *optimism in the face of uncertainty* (OFU) principle

Background – contextual linear bandit

$$r_{t,a_t} = h(\mathbf{x}_{t,a_t}) + \epsilon_t, \quad \epsilon_t \text{ -sub-Gaussian}$$

- | Build confidence set for θ and use *optimism in the face of uncertainty* (OFU) principle
- | Leads to $\Theta(d^{\rho} \bar{T})$ regret (Abbasi-Yadkori et al. 2011)
- | Strongly depends on linear structure!

Background – general reward function

$$r_{t;a_t} = h(\mathbf{x}_{t;a_t}) + \epsilon_t, \quad 0 \leq h(\mathbf{x}) \leq 1, \quad \epsilon_t \text{ -sub-Gaussian}$$

Background – general reward function

$$r_{t;a_t} = h(\mathbf{x}_{t;a_t}) + \epsilon_t, \quad 0 \leq h(\mathbf{x}) \leq 1, \quad \epsilon_t \text{ -sub-Gaussian}$$

| Including many popular contextual bandit problems

| Linear bandit

| $h(\mathbf{x}) = \mathbf{w}^\top \mathbf{x}$, where $\|\mathbf{w}\|_2 \leq 1, \|\mathbf{x}\|_2 \leq 1$

| Generalized linear bandit

| $h(\mathbf{x}) = g(\mathbf{w}^\top \mathbf{x})$, where $\|\mathbf{w}\|_2 \leq 1, \|\mathbf{x}\|_2 \leq 1, |g| \leq 1$

Background – general reward function

$$r_{t;a_t} = h(\mathbf{x}_{t;a_t}) + \epsilon_t, \quad 0 \leq h(\mathbf{x}) \leq 1, \quad \epsilon_t \text{ -sub-Gaussian}$$

| Including many popular contextual bandit problems

| Linear bandit

| $h(\mathbf{x}) = \mathbf{w}^\top \mathbf{x}$, where $\|\mathbf{w}\|_2 \leq 1, \|\mathbf{x}\|_2 \leq 1$

| Generalized linear bandit

| $h(\mathbf{x}) = g(\mathbf{w}^\top \mathbf{x})$, where $\|\mathbf{w}\|_2 \leq 1, \|\mathbf{x}\|_2 \leq 1, |g| \leq 1$

We do not know what h is...

Background – general reward function

$$r_{t;a_t} = h(\mathbf{x}_{t;a_t}) + \epsilon_t, \quad 0 \leq h(\mathbf{x}) \leq 1, \quad \epsilon_t \text{ -sub-Gaussian}$$

| Including many popular contextual bandit problems

| Linear bandit

$$| \quad h(\mathbf{x}) = \mathbf{w}^\top \mathbf{x}, \text{ where } \|\mathbf{w}\|_2 \leq 1, \|\mathbf{x}\|_2 \leq 1$$

| Generalized linear bandit

$$| \quad h(\mathbf{x}) = g(\mathbf{w}^\top \mathbf{x}), \text{ where } \|\mathbf{w}\|_2 \leq 1, \|\mathbf{x}\|_2 \leq 1, |g| \leq 1$$

We do not know what h is...

Use some universal function approximator, such as neural networks!

Background – neural network

Fully connected neural networks:

$$f(\mathbf{x}; \theta) = \rho\left(\sum_{l=1}^L \mathbf{W}_l \mathbf{x}\right)$$

Background – neural network

Fully connected neural networks:

$$f(\mathbf{x}; \theta) = \rho \left(\sum_{l=1}^L \mathbf{W}_l \mathbf{x} \right)$$

| $\rho(x) = \max\{x, 0\}$ is the ReLU activation function

Background – neural network

Fully connected neural networks:

$$f(\mathbf{x}; \theta) = \rho_{\overline{m}} \mathbf{W}_L \mathbf{W}_{L-1} \dots \mathbf{W}_1 \mathbf{x}$$

- | $\rho(x) = \max\{x, 0\}$ is the ReLU activation function
- | \mathbf{W}_i is the weight matrix
 - | $\mathbf{W}_1 \in \mathbb{R}^{m \times d}$
 - | $\mathbf{W}_i \in \mathbb{R}^{m_i \times m_{i-1}}$ $i = 2, \dots, L-1$
 - | $\mathbf{W}_L \in \mathbb{R}^{m \times 1}$

Background – neural network

Fully connected neural networks:

$$f(\mathbf{x}; \theta) = \rho_m \left(\mathbf{W}_L \mathbf{W}_{L-1} \dots \mathbf{W}_1 \mathbf{x} \right)$$

- | $\rho(x) = \max\{x, 0\}$ is the ReLU activation function
- | $\theta = [\text{vec}(\mathbf{W}_1)^T; \dots; \text{vec}(\mathbf{W}_L)^T]^T \in \mathbb{R}^p$, $p = m + md + m^2(L-1)$

Background – neural network

Fully connected neural networks:

$$f(\mathbf{x}; \theta) = \rho_m \mathbf{W}_L \mathbf{W}_{L-1} \dots \mathbf{W}_1 \mathbf{x}$$

- | $\rho(x) = \max\{x, 0\}$ is the ReLU activation function
- | $\theta = [\text{vec}(\mathbf{W}_1)^T; \dots; \text{vec}(\mathbf{W}_L)^T]^T \in \mathbb{R}^p$, $p = m + md + m^2(L-1)$
- | Gradient of the neural network $\mathbf{g}(\mathbf{x}; \theta) = \nabla_{\theta} f(\mathbf{x}; \theta) \in \mathbb{R}^p$

Question

- | Neural network-based contextual bandit algorithms (Riquelme et al. 2018; Zahavy and Mannor 2019)
- | No theoretical guarantee

Question

- | Neural network-based contextual bandit algorithms (Riquelme et al. 2018; Zahavy and Mannor 2019)
- | No theoretical guarantee

Can we design provably efficient neural network-based algorithm to learn the general reward function?

Question

- | Neural network-based contextual bandit algorithms (Riquelme et al. 2018; Zahavy and Mannor 2019)
- | No theoretical guarantee

Can we design provably efficient neural network-based algorithm to learn the general reward function?

Yes! NeuralUCB

- | Neural network to model reward function, UCB strategy to explore
- | Theoretical guarantee on regret $\mathcal{O}(\sqrt{pT})$
- | Matches regret bound for linear setting (Abbasi-Yadkori et al. 2011)

NeuralUCB { initialization

| Special initialization on 0

| For $1 \leq l \leq L - 1$,

$$W_l = \begin{pmatrix} W & 0 \\ 0 & W \end{pmatrix}; W_{f_{i,j}^g} \sim N(0; 4/m)$$

| For L , $W = (w^>; w^>)$, $w_{f_{i,j}^g} \sim N(0; 2/m)$

NeuralUCB { initialization

| Special initialization on 0

| For $1 \leq l \leq L-1$,

$$W_l = \begin{pmatrix} W & 0 \\ 0 & W \end{pmatrix}; W_{f(i,j)} \sim N(0; 4^{-l}m)$$

| For L , $W = \begin{pmatrix} w^> & \\ & w^> \end{pmatrix}$, $w_{f(i,j)} \sim N(0; 2^{-l}m)$

| Normalization on x^i 's: for any $1 \leq i \leq TK$, $\|x^i\|_2 = 1$ and $[x^i]_j = [x^i]_{j+d=2}$

| For any unit vector x , construct $x^0 = \begin{pmatrix} x \\ x \end{pmatrix} = \frac{1}{\sqrt{2}}$

NeuralUCB { initialization

- Special initialization on w^0

- For $1 \leq l \leq L-1$,

$$W_l = \begin{pmatrix} W & 0 \\ 0 & W \end{pmatrix}; W_{f(i,j)} \sim N(0; 4\sigma^2 m)$$

- For L , $W = \begin{pmatrix} w^> & \\ & w^> \end{pmatrix}$, $w_{f(i,j)} \sim N(0; 2\sigma^2 m)$

- Normalization on x^i 's: for any $1 \leq i \leq TK$, $\|x^i\|_2 = 1$ and $[x^i]_j = [x^i]_{j+d=2}$

- For any unit vector x , construct $x^0 = \begin{pmatrix} x \\ x \end{pmatrix} = \frac{1}{\sqrt{2}}$

Guarantee that $\langle x^i; x^0 \rangle = 0$!

NeuralUCB { upper confidence bounds

At round t , NeuralUCB will...

- | Observe $x_{t;a} g_{a=1}^K$

NeuralUCB { upper con dence bounds

At round t , NeuralUCB will...

- | Observe $x_{t;a} g_{a=1}^K$
- | Compute upper con dence bound for each a , which is

$$U_{t;a} = \underbrace{\bar{g}(x_{t;a}; t-1)}_{\text{mean}} + \underbrace{t^{-1} \sqrt{q \frac{g(x_{t;a}; t-1) > Z_t^{-1} g(x_{t;a}; t-1) = m}{\text{variance}}}}_{\text{variance}}$$

NeuralUCB { upper con dence bounds

At round t , NeuralUCB will...

- | Observe $x_{t;a} g_{a=1}^K$
- | Compute upper con dence bound for each arm, which is

$$U_{t;a} = \underbrace{\left\{ \frac{1}{Z_t} \sum_{j=1}^t x_{t;a;j} \right\}}_{\text{mean}} + \underbrace{t^{-1} \sqrt{q \frac{g(x_{t;a}; t-1) > Z_t^{-1} \sum_{j=1}^t g(x_{t;a;j})}{Z_t^{-1} \sum_{j=1}^t g(x_{t;a;j})^2}}}_{\text{variance}}$$

Compared with LinUCB (Li et al. 2010)

$$U_{t;a} = \underbrace{\left\{ \frac{1}{Z_t} \sum_{j=1}^t x_{t;a;j} \right\}}_{\text{mean}} + \underbrace{t^{-1} \sqrt{q \frac{x_{t;a}^T Z_t^{-1} x_{t;a}}{Z_t^{-1} \sum_{j=1}^t x_{t;a;j}^2}}}_{\text{variance}}$$

NeuralUCB { upper con dence bounds

At round t , NeuralUCB will...

- Observe $x_{t;a} g_{a=1}^K$
- Compute upper con dence bound for each arm, which is

$$U_{t;a} = \underbrace{\left[\frac{1}{Z_{t-1}} \sum_{j=1}^t x_{t;a;j} \right]}_{\text{mean}} + \underbrace{t-1 \sqrt{\frac{q}{Z_{t-1} \sum_{j=1}^t g(x_{t;a;j})^2}}}_{\text{variance}}$$

Compared with LinUCB (Li et al. 2010)

$$U_{t;a} = \underbrace{\left[\frac{1}{Z_{t-1}} \sum_{j=1}^t x_{t;a;j} \right]}_{\text{mean}} + \underbrace{t-1 \sqrt{\frac{q}{Z_{t-1} \sum_{j=1}^t x_{t;a;j}^2}}}_{\text{variance}}$$

- Select $a_t = \operatorname{argmax}_{a \in [K]} U_{t;a}$, play a_t and observe reward $r_{t;a_t}$

NeuralUCB { update parameter

After receiving reward, NeuralUCB will...

- | Update Z_t

$$Z_t = Z_{t-1} + g(x_{t;a_t}; \theta_{t-1})g(x_{t;a_t}; \theta_{t-1})^{\top} = m$$

NeuralUCB { update parameter

After receiving reward, NeuralUCB will...

- | Update Z_t

$$Z_t = Z_{t-1} + g(x_{t;a_t}; \theta_{t-1})g(x_{t;a_t}; \theta_{t-1})^\top = m$$

- | Update θ_t using gradient descent

- | Denote loss function $L(\theta)$ as

$$L(\theta) = \sum_{i=1}^X (f(x_{i;a_i}; \theta) - r_{i;a_i})^2 = 2 + m k \quad (0) k_2^2 = 2$$

NeuralUCB { update parameter

After receiving reward, NeuralUCB will...

- | Update Z_t

$$Z_t = Z_{t-1} + g(x_{t;a_t}; \theta_{t-1}) g(x_{t;a_t}; \theta_{t-1})^T = m$$

- | Update θ_t using gradient descent

- | Denote loss function $L(\theta)$ as

$$L(\theta) = \sum_{i=1}^X (f(x_{i;a_i}; \theta) - r_{i;a_i})^2 = 2 + m k \quad (0) k_2^2 = 2$$

- | Run J step gradient descent on $L(\theta)$ starting from θ_0 , take θ_t as the last iterate

$$\theta^{(0)} = \theta_0; \quad \theta^{(j+1)} = \theta^{(j)} - \eta \nabla L(\theta^{(j)}); \quad \theta_t = \theta^{(J)}$$

NeuralUCB { confidence radius

After update neural network function, NeuralUCB will compute, which is ...

Under the overparameterized setting ($m \gg 1$),

$$t = O \left(\underbrace{p}_{S} + \underbrace{r \log \frac{\det Z_t}{\det}}_{\text{confidence radius}} + \underbrace{\left(\frac{+ tL}{\{z^m\}} \right)^{J=2}}_{\text{function approximation error}} \right)$$

NeuralUCB { confidence radius

After update neural network function, NeuralUCB will compute, which is ...

Under the overparameterized setting ($m \gg 1$),

$$t = O \left(\underbrace{p_S + \frac{r}{\log \frac{\det Z_t}{\det I}}}_{\text{confidence radius}} + \underbrace{\left(\frac{+ tL}{\det I} \right)^{J=2} \frac{p}{t}}_{\text{function approximation error}} \right)$$

Compared with LinUCB,

$$t = O \left(p_S + \frac{r}{\log \frac{\det Z_t}{\det I}} \right)$$

no function approximation error part!

Main theory { assumptions

Assumption

There exists $\epsilon_0 > 0$ such that $\|H - H_0\| \leq \epsilon_0$, where H is the neural tangent kernel matrix (Jacot et al. 2018; Cao and Gu 2019) on contexts $\{x^i\}_{i=1}^n$.

Main theory { assumptions

Assumption

There exists $\epsilon_0 > 0$ such that $\|H - \epsilon_0 I\| \leq \epsilon_0$, where H is the neural tangent kernel matrix (Jacot et al. 2018; Cao and Gu 2019) on contexts $\{x^i\}_{i=1}^K$.

- | Satisfied if no two contexts $\{x^i\}_{i=1}^K$ are parallel.

Main theory { assumptions

Assumption

There exists $\epsilon_0 > 0$ such that $\|H - \epsilon_0 I\| \leq \epsilon_0$, where H is the neural tangent kernel matrix (Jacot et al. 2018; Cao and Gu 2019) on contexts $\{x^i\}_{i=1}^n$.

| Satisfied if no two contexts $\{x^i\}_{i=1}^n$ are parallel.

Definition

The effective dimension d_{eff} of the neural tangent kernel matrix on contexts $\{x^i\}_{i=1}^n$ is defined as $d_{\text{eff}} = \log \det(I + H) = \log(1 + \text{Tr}(H))$.

Main theory { assumptions

Assumption

There exists $\epsilon_0 > 0$ such that $\|H\| \leq \epsilon_0$, where H is the neural tangent kernel matrix (Jacot et al. 2018; Cao and Gu 2019) on contexts $\{x^i\}_{i=1}^n$.

- | Satisfied if no two contexts $\{x^i\}_{i=1}^n$ are parallel.

Definition

The effective dimension d_{eff} of the neural tangent kernel matrix on contexts $\{x^i\}_{i=1}^n$ is defined as $d_{\text{eff}} = \log \det(I + H) = \log(1 + \text{TK})$.

- | Notion adapted from Valko et al. (2013) and Yang and Wang (2019)

Main theory { assumptions

Assumption

There exists $\epsilon_0 > 0$ such that $\|H\| \leq \epsilon_0$, where H is the neural tangent kernel matrix (Jacot et al. 2018; Cao and Gu 2019) on contexts $\{x^i\}_{i=1}^T$.

- | Satisfied if no two contexts $\{x^i\}_{i=1}^T$ are parallel.

Definition

The effective dimension \mathfrak{d} of the neural tangent kernel matrix on contexts $\{x^i\}_{i=1}^T$ is defined as $\mathfrak{d} = \log \det(I + H) = \log(1 + \text{TK})$.

- | Notion adapted from Valko et al. (2013) and Yang and Wang (2019)
- | $\mathfrak{d} \approx \log T$ in several special cases (Valko et al. 2013)

Main theory { regret bound

Theorem

Let $h = [h(x^i)]_{i=1}^{TK} \in \mathbb{R}^{TK}$. Set $J = \frac{e}{2} (TL + \frac{1}{m})$,
 $S = \frac{1}{2} \frac{e}{h^2 H^2}$. Under the
 overparameterized setting $(TL > 1)$, with probability at least $1 - \frac{1}{m}$,

$$R_T = \mathcal{O} \left(\frac{1}{\sqrt{TL}} \max_{\theta \in \mathcal{S}} \|\theta\|_2 \right)$$

Main theory – regret bound

Theorem

Let $\mathbf{h} = [h(\mathbf{x}^i)]_{i=1}^{TK} \in \mathbb{R}^{TK}$. Set $J = \frac{e}{2} (TL + m)$,
 $= ((mTL + m) \cdot 1)$ and $S = \frac{1}{2} \frac{e}{\mathbf{h}^T \mathbf{H}^{-1} \mathbf{h}}$. Under the
 overparameterized setting ($m \gg 1$), with probability at least $1 - \delta$,

$$R_T = \mathcal{O} \left(\frac{1}{\sqrt{T}} \max_{\theta} \int S^2 g \right)$$

| h belongs to the RKHS space H spanned by \mathbf{H}) $S \propto \frac{1}{\|\mathbf{h}\|_H}$

Main theory – regret bound

Theorem

Let $\mathbf{h} = [h(\mathbf{x}^i)]_{i=1}^{TK} \in \mathbb{R}^{TK}$. Set $J = \frac{e}{\rho} (TL + m)$,
 $S = \frac{1}{\sqrt{mTL + m}}$ and $S = \frac{1}{\sqrt{\mathbf{h}^T \mathbf{H}^{-1} \mathbf{h}}}$. Under the
 overparameterized setting ($m \gg 1$), with probability at least $1 - \delta$,

$$R_T = \mathcal{O} \left(\frac{\rho}{\sqrt{mTL + m}} \max_{\theta} \int_0^1 \langle \theta; S^2 g \rangle \right)$$

- | h belongs to the RKHS space H spanned by \mathbf{H}) $S = \frac{1}{\sqrt{\mathbf{h}^T \mathbf{H}^{-1} \mathbf{h}}}$
- | R_T does not depend on ρ , the dimension of the dynamic feature mapping $\mathbf{g}(\mathbf{x}; \theta)$

Main theory – regret bound

Theorem

Let $\mathbf{h} = [h(\mathbf{x}^i)]_{i=1}^{TK} \in \mathbb{R}^{TK}$. Set $J = \frac{e}{2} (TL + m)$,
 $= ((mTL + m) \cdot 1)$ and $S = \frac{2}{\mathbf{h}^T \mathbf{H}^{-1} \mathbf{h}}$. Under the
 overparameterized setting ($m \gg 1$), with probability at least $1 - \delta$,

$$R_T = \Theta \left(\frac{\rho}{\delta T} \max_{\theta} \int S^2 g \right)$$

- | h belongs to the RKHS space H spanned by \mathbf{H} $\|h\|_H$
- | R_T does not depend on ρ , the dimension of the dynamic feature mapping $\mathbf{g}(\mathbf{x}; \cdot)$
- | Recover the regret for linear contextual bandit $\Theta(d^{\rho} \bar{T})$
 (Abbasi-Yadkori et al. 2011)

Takeaway message

- | NeuralUCB uses neural network $f(\mathbf{x}; t)$ to predict, gradient $\mathbf{g}(\mathbf{x}; t)$ to explore

Takeaway message

- | NeuralUCB uses neural network $f(\mathbf{x}; t)$ to predict, gradient $\mathbf{g}(\mathbf{x}; t)$ to explore
- | NeuralUCB achieves $\Theta(\sqrt{T})$ regret, matches result for linear setting

Takeaway message

- | NeuralUCB uses neural network $f(\mathbf{x}; t)$ to predict, gradient $\mathbf{g}(\mathbf{x}; t)$ to explore
- | NeuralUCB achieves $\Theta(\sqrt{T})$ regret, matches result for linear setting
- | NeuralUCB also works well empirically

Takeaway message

- | NeuralUCB uses neural network $f(\mathbf{x}; t)$ to predict, gradient $\mathbf{g}(\mathbf{x}; t)$ to explore
- | NeuralUCB achieves $\Theta(\sqrt{T})$ regret, matches result for linear setting
- | NeuralUCB also works well empirically

Thank you!