

به نام خدا

گزارش پروژه درس یادگیری ماشین

استاد درس:

دکتر جواد سلیمی

دانشجو:

مهرداد چهارباغی

دانشگاه کاشان

دانشکده برق و کامپیوتر

زمستان 1402

## پروژه اول:

الگوریتم‌های Contextual bandit و پیاده‌سازی آن‌ها را با Pearl و با استفاده از یک محیط Contextual bandit برای ارائه داده‌ها از مجموعه داده‌های UCI نشان می‌دهیم و عملکرد SquareCB، LinUCB و LinTS را آزمایش می‌کنیم.

ابتدا باید توضیحاتی راجع به این سه و همچنین خود Contextual bandit داده شود که در ادامه آورده شده است.

### • Contextual bandit:

نوعی مسئله در یادگیری تقویتی است. که در این الگوریتم نه تنها باید از بین اکشن‌های مختلف (یا "Bandit") برای به حداکثر رساندن پاداش خود انتخاب کند، بلکه contextual (یا حالت) را نیز در نظر می‌گیرد که در آن تصمیم گرفته می‌شود. بخش "contextual" به اطلاعات اضافی موجود در زمان تصمیم‌گیری اشاره دارد که می‌تواند برای پیش‌بینی نتیجه یا پاداش هر اقدام با دقت بیشتری مورد استفاده قرار گیرد.

در این مسئله، قبل از تصمیم‌گیری، الگوریتم با اطلاعات Contextual مربوط به هر تصمیم ارائه می‌شود. این Context می‌تواند به طور قابل توجهی بر پاداش مورد انتظار یک عمل تأثیر بگذارد.

چالش اصلی در مشکلات Contextual bandit ایجاد تعادل بین اکتشاف<sup>1</sup> (آزمودن اقدامات مختلف برای جمع‌آوری اطلاعات بیشتر در مورد پاداش‌های آنها) و بهره‌برداری<sup>2</sup> (انتخاب شناخته‌شده‌ترین اقدام برای به حداکثر رساندن پاداش) است. الگوریتم باید از تصمیمات گذشته و نتایج آنها درس بگیرد و استراتژی خود را بر اساس بازخورد پاداش و Context ای که در آن تصمیم گرفته شده است به روز کند.

### • SquareCB

الگوریتم SquareCB فقط به یک مدل رگرسیون نیاز دارد که با آن تابع پاداش<sup>3</sup> را یاد می‌گیرد. با توجه به مدل پاداش، SquareCB سیاست زیر را اجرا می‌کند.

---

<sup>1</sup> Exploration

<sup>2</sup> Exploitation

<sup>3</sup> Reward

$$\hat{a}_* \in \arg \max_a \hat{r}(x, a)$$

$$\hat{r}_* \in \max_a \hat{r}(x, a)$$

$$\text{If } a \neq \hat{a}_* : \pi(a, x) = \frac{1}{A + \gamma(\hat{r}_* - \hat{r}(x, a))}$$

$$\text{If } a = \hat{a}_* : \pi(a, x) = 1 - \sum_{a' \neq \hat{a}_*} \pi(a', x).$$

این سیاست اکتشافی که اکتشاف<sup>4</sup> و بهره برداری<sup>5</sup> را به شیوه ای هوشمندانه متعادل می کند.

برای استفاده از الگوریتم SquareCB در Pearl، یادگیرنده سیاست<sup>6</sup> را به عنوان NeuralBandit تنظیم می کنیم. NeuralBandit یک کلاس پایه است و از تخمین تابع پاداش با معماری عصبی پشتیبانی می کند. با داشتن دسترسی به یک مدل پاداش تخمینی، می توانیم مازول کاوش را با مازول SquareCBExploration نمونه سازی می کنیم.

## • LinUCB

در Pearl از نسخه عصبی الگوریتم LinUCB استفاده می شود که از نوع UCB کاوش با معماری های عصبی استفاده می کند. LinUCB و نسخه عصبی آن، تعمیم هایی از الگوریتم UCB (Upper Confidence Bound) هستند. هر دو یک سیاست به شکل زیر را اجرا می کنند:

$$\pi(a, x) \in \arg \max_a \hat{r}(x, a) + \text{score}(x, a)$$

یعنی، هر دو از تابعی استفاده می کنند که پاداش مورد انتظار را با یک عبارت پاداش اضافی تخمین می زند، که پتانسیل انتخاب یک عمل را با توجه به زمینه خاصی کمیت می دهد. یک روش معمول برای تخمین تابع امتیاز، در حالت خطی، زمانی که ویژگی ها  $\phi(x, a)$  هستند، از طریق:

$$\text{score}(x, a) = \alpha \|\phi(x, a)\|_{A^{-1}}$$

$$A = \lambda I + \sum_{n \leq t} \phi(x_n, a_n) \phi^T(x_n, a_n) \quad \text{به طوریکه}$$

<sup>4</sup> Exploration

<sup>5</sup> Exploitation

<sup>6</sup> Policy Learner

## • LinTS

در ادامه، نحوه استفاده از نسخه عصبی الگوریتم LinTS را با کتابخانه Pearl، یعنی الگوریتمی که از کاوش نمونه‌برداری تامپسون<sup>7</sup> با معماری‌های عصبی استفاده می‌کند، توضیح می‌دهیم. نمونه برداری LinTS ارتباط نزدیکی با الگوریتم LinUCB دارد، با یک اصلاح کلیدی که اغلب همگرایی آن را در عمل بهبود می‌بخشد: *تابع امتیاز را از یک احتمال نمونه برداری کنید، به جای اینکه آن را به طور قطعی ثابت کنید.* در عمل، اینکار اغلب بازوهای کاوش بیش از حد را کاهش می‌دهد، زیرا امتیاز ممکن است کمتر از الگوریتم LinUCB باشد.

### مراحل:

#### • نصب:

```
• %pip uninstall Pearl -y
• %rm -rf Pearl
• !git clone https://github.com/facebookresearch/Pearl.git
• %cd Pearl
• %pip install .
• %cd ..
```

ابتدا باید با دستوراتی که در کد بالا آورده شده، پکیج Pearl را نصب کنیم.

#### • وارد کردن ماژول‌ها و کتابخانه‌ها:

چون تعداد ماژول‌هایی که باید وارد شوند زیاد است از نوشتن آنها اینجا خودداری کرده چون تمام آنها در کد پیاده‌سازی آورده شده است.

#### • بارگذاری محیط:

```
1. # load environment
2. device = torch.device("cuda" if torch.cuda.is_available() else
   "cpu")
3.
4. # Download UCI dataset
5. uci_data_path = "./utils/instantiations/environments/uci_datasets"
6. if not os.path.exists(uci_data_path):
7.     os.makedirs(uci_data_path)
8.     download_uci_data(data_path=uci_data_path)
```

ابتدا باید دیتاست مورد نیاز را دریافت کرده تا بتوانیم از آن استفاده کنیم. این کار در خطوط 4 تا 8 انجام شده است.

---

<sup>7</sup> Thompson sampling exploration

```
Built CB environment using the pendigits UCI dataset
```

```
pendigits_uci_dict = {  
    "path_filename": os.path.join(uci_data_path,  
    "pendigits/pendigits.tra"),  
    "action_embeddings": "discrete",  
    "delim_whitespace": False,  
    "ind_to_drop": [],  
    "target_column": 16,  
}
```

```
env = SLCBEnvironment(**pendigits_uci_dict)
```

خب نوبت ساخت محیط Contextual Bandit می‌رسد. در اینجا از دیتاست pendigit استفاده شده است.

```
number_of_steps = 10000
```

```
record_period = 400
```

و در اینجا هم تعداد قدم‌ها (مراحل).

### • پیاده سازی SquareCB

با توجه به طراحی ماژولار Pearl، ما از OneHotActionTensorRepresentationModule به عنوان ماژول نمایش اکشن استفاده می‌کنیم. یعنی، زمانی که مجموعه عمل تعداد محدودی از عناصر  $\{1, 2, \dots, N\}$  به عنوان یک بردار One-hot باشد.

```
1. action_representation_module =  
    OneHotActionTensorRepresentationModule(  
2.     max_number_actions= env._action_space.n,  
3. )  
4.  
5. agent = PearlAgent(  
6.     policy_learner=NeuralBandit(  
7.         feature_dim = env.observation_dim + env._action_space.n,  
8.         hidden_dims=[64, 16],  
9.         training_rounds=50,  
10.         action_representation_module=action_representation_mod  
    ule,  
11.         exploration_module= SquareCBExploration(gamma =  
        env.observation_dim * env._action_space.n * number_of_steps)  
12.     ),  
13.     replay_buffer=FIFOOffPolicyReplayBuffer(100_000),  
14.     device_id=-1,  
15. )
```

در خط اول نوع اکشن و فضای آن مشخص شده است.

در خط 5 تا 12 : ویژگی ایجنت خود را مشخص میکنیم.در اینجا سیاست یادگیری ما از نوع عصبی خواهد بود.  
و تابع اکتشاف هم از نوع SquareCB.

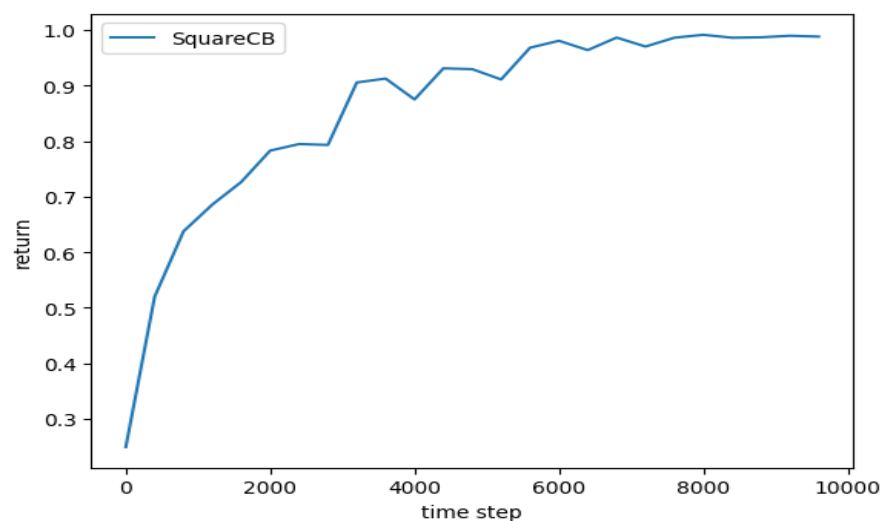
در ادامه:

```
info = online_learning(  
    agent=agent,  
    env=env,  
    number_of_steps=number_of_steps,  
    print_every_x_steps=100,  
    record_period=record_period,  
    learn_after_episode=True,  
)
```

همه ویژگی ها رو یکجا جمع می کنیم.

و اما در آخر نتیجه را نمایش می دهیم:

```
torch.save(info["return"], "SquareCB-return.pt")  
plt.plot(record_period * np.arange(len(info["return"])), info["return"],  
label="SquareCB")  
plt.xlabel("time step")  
plt.ylabel("return")  
plt.legend()  
plt.show()
```



## LinUCB •

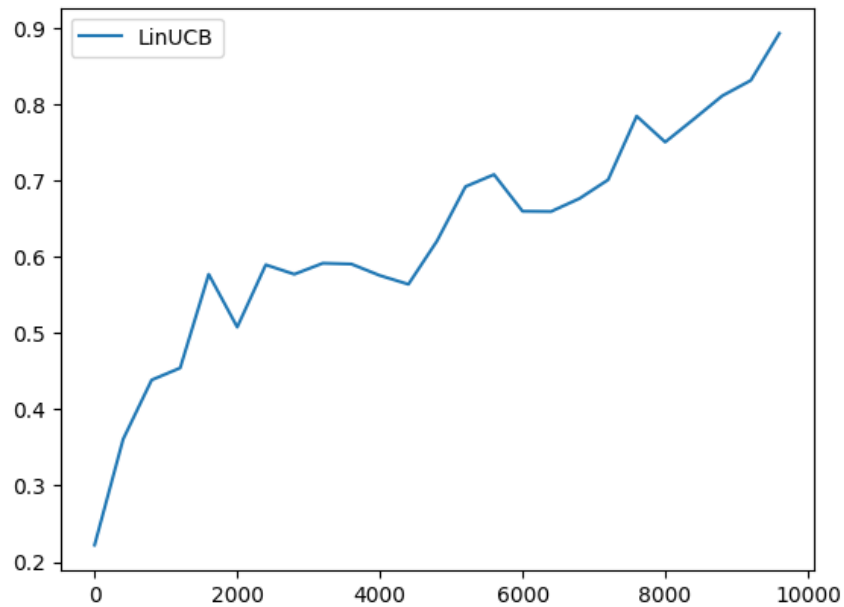
برای پیاده سازی الگوریتم LinUCB در Pearl، از ماژول یادگیرنده با سیاست **NeuralLinearBandit** استفاده می کنیم. این ماژول از (i) یادگیری یک مدل پاداش پشتیبانی می کند، و (ii) یک تابع امتیاز را با تخمین عدم قطعیت با استفاده از ویژگی های لایه آخر محاسبه می کند. علاوه بر این، ماژول کاوش را روی **UCBExploration** تنظیم می کنیم.

```
agent = PearlAgent(  
    policy_learner=NeuralLinearBandit(  
        feature_dim = env.observation_dim + env._action_space.n,  
        hidden_dims=[64, 16],  
        training_rounds=50,  
        action_representation_module=action_representation_module,  
        exploration_module= UCBExploration(alpha=1.0)  
    ),  
    replay_buffer=FIFOOffPolicyReplayBuffer(100_000),  
    device_id=-1,  
)  
info = online_learning(  
    agent=agent,  
    env=env,  
    number_of_steps=number_of_steps,  
    print_every_x_steps=100,  
    record_period=record_period,  
    learn_after_episode=True,  
)
```

در اینجا هم نوع سیاست یادگیر را روی عصبی خطی تنظیم می کنیم و ماژول اکتشاف را روی UCB.

و نمایش نتیجه:

```
torch.save(info["return"], "LinUCB-return.pt")  
plt.plot(record_period * np.arange(len(info["return"])), info["return"],  
label="LinUCB")  
plt.xlabel("time step")  
plt.ylabel("return")  
plt.legend()  
plt.show()
```



## • LinTS

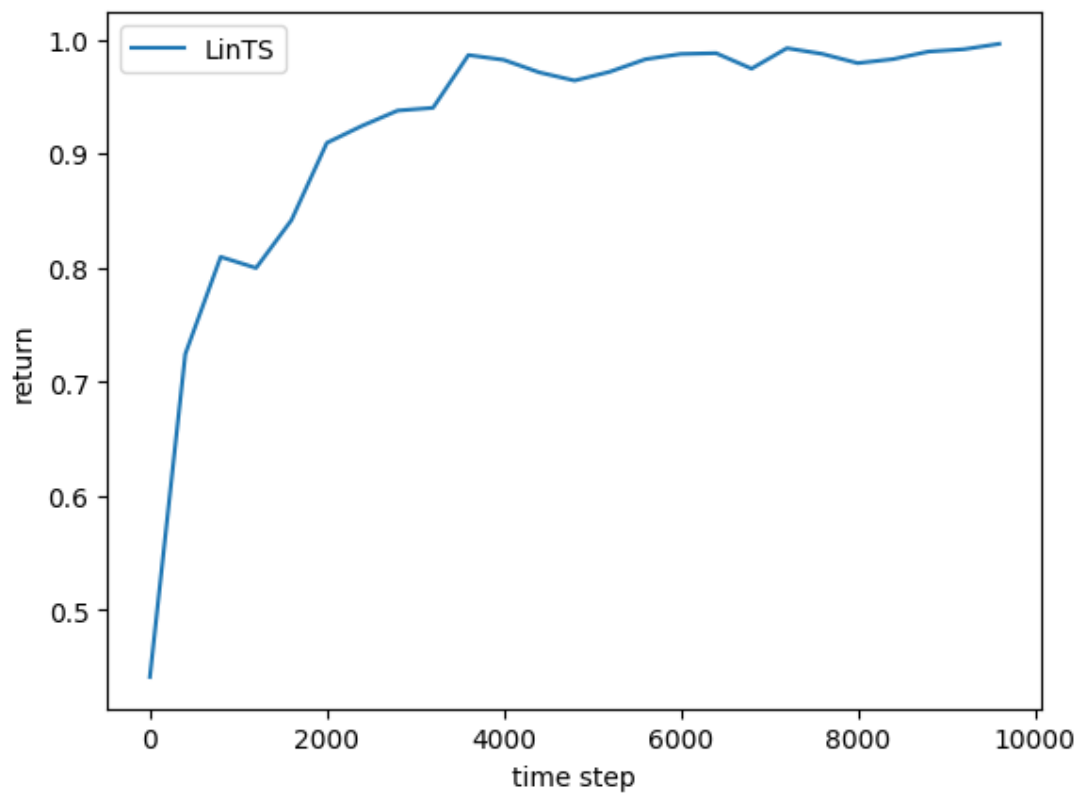
برای پیاده سازی الگوریتم LinTS در Pearl، از ماژول یادگیری سیاست NeuralLinearBandit استفاده می-کنیم. علاوه بر این، ماژول کاوش را روی ThompsonSamplingExplorationLinear تنظیم می کنیم. این عامل را قادر می سازد تا امتیاز را بر اساس عدم قطعیت تخمینی خود نمونه برداری کند، به جای اینکه آن را مانند الگوریتم LinUCB تصحیح کند.

```
action_representation_module = OneHotActionTensorRepresentationModule(
    max_number_actions= env._action_space.n,
)

agent = PearlAgent(
    policy_learner=NeuralLinearBandit(
        feature_dim = env.observation_dim + env._action_space.n,
        hidden_dims=[64, 16],
        training_rounds=50,
        action_representation_module=action_representation_module,
        exploration_module= ThompsonSamplingExplorationLinear()
    ),
    replay_buffer=FIFOOffPolicyReplayBuffer(100_000),
    device_id=-1,
)
```



```
info = online_learning(  
    agent=agent,  
    env=env,  
    number_of_steps=number_of_steps,  
    print_every_x_steps=100,  
    record_period=record_period,  
    learn_after_episode=True,  
)  
torch.save(info["return"], "LinTS-return.pt")  
plt.plot(record_period * np.arange(len(info["return"])), info["return"],  
label="LinTS")  
plt.xlabel("time step")  
plt.ylabel("return")  
plt.legend()  
plt.show()
```



## پروژه دوم:

### پیاده سازی دریاچه یخ زده<sup>8</sup> با استفاده از Pearl

این پروژه نحوه استفاده از DQN را برای حل محیط FrozenLake-v1 نشان می دهد. این محیط مشاهداتی<sup>9</sup> را به عنوان دارد. در ادامه، نحوه استفاده از پکیج و ماژول Pearl's OneHotObservationsFromDiscrete را برای تبدیل مشاهدات به نمایش های One-hot آنها نشان می دهیم. در ادامه همراه با کد پیاده سازی قدم به قدم جلو می رویم و توضیح داده می شود که هر قسمت چه کاری انجام می دهد.

#### مراحل:

##### • نصب:

```
• %pip uninstall Pearl -y
• %rm -rf Pearl
• !git clone https://github.com/facebookresearch/Pearl.git
• %cd Pearl
• %pip install .
• %cd ..
```

ابتدا باید با دستوراتی که در کد بالا آورده شده، پکیج Pearl را نصب کنیم.

##### • وارد کردن ماژول ها و کتابخانه های مورد نیاز:

```
• from pearl.utils.functional_utils.experimentation.set_seed import set_seed
• from pearl.policy_learners.sequential_decision_making.deep_q_learning import
  DeepQLearning
• from pearl.replay_buffers.sequential_decision_making.fifo_off_policy_replay_buffer import
  FIFOffPolicyReplayBuffer
• from pearl.utils.functional_utils.train_and_eval.online_learning import online_learning
• from pearl.pearl_agent import PearlAgent
• from pearl.utils.instantiations.environments.gym_environment import GymEnvironment
• from pearl.utils.instantiations.environments.environments import (
•     OneHotObservationsFromDiscrete,
• )
• from pearl.utils.instantiations.spaces.discrete import DiscreteSpace
• import torch
```

---

<sup>8</sup> Frozen Lake

<sup>9</sup> Observation

- `import matplotlib.pyplot as plt`
- `import numpy as np`
- `from pearl.action_representation_modules.one_hot_action_representation_module`
- `import (OneHotActionTensorRepresentationModule)`
- `set_seed(0)`

• حل مسئله:

مسئله را با استفاده از Vanilla DQN حل می‌کنیم:

```
number_of_steps = 20000
record_period = 400
```

ابتدا تعداد قدم<sup>10</sup> هارو مشخص می‌کنیم.

```
env = OneHotObservationsFromDiscrete(
    GymEnvironment(
        "FrozenLake-v1", is_slippery=False, map_name="4x4",
    )
)
```

در اینجا مسئله ای را که می‌خواهیم حل کنیم ویژگی‌های آن را مشخص می‌کنیم. (به طور مثال ابعاد زمین، لغزنده بودن...)

```
action_representation_module = OneHotActionTensorRepresentationModule(
    max_number_actions= env.action_space.n,
)
```

تعیین نوع اکشن و فضای آن

```
1 state_dim = env.observation_space.n
2 agent = PearlAgent(
3     policy_learner=DeepQLearning(
```

---

<sup>10</sup> Step

```

4     state_dim=state_dim,
5     action_space=env.action_space,
6     hidden_dims=[64, 64],
7     training_rounds=1,
8     action_representation_module=action_representation_module
9 ),
10    replay_buffer=FIFOOffPolicyReplayBuffer(1000),
11 )

```

در اینجا در خط اول فضای ابعاد حالت را مشخص می‌کنیم.

در خط 2 تا 9 نوع ایجنت<sup>11</sup> که در اینجا از نوع ایجن Pearl قاعدتا باید باشد به همراه سیاستی که باید پیش بگیرد و دیگر ویژگی‌ها آن را مشخص می‌کنیم.

```

info = online_learning(
    agent=agent,
    env=env,
    number_of_steps=number_of_steps,
    print_every_x_steps=100,
    record_period=record_period,
    learn_after_episode=False,
)

```

در اینجا نوع یادگیری که اینجا آنلاین هست و آرگومانهای ورودی آن مثل ایجنت که از قبل مشخص کردیم، محیط، تعداد قدم‌ها و ... مشخص می‌کنیم.

```

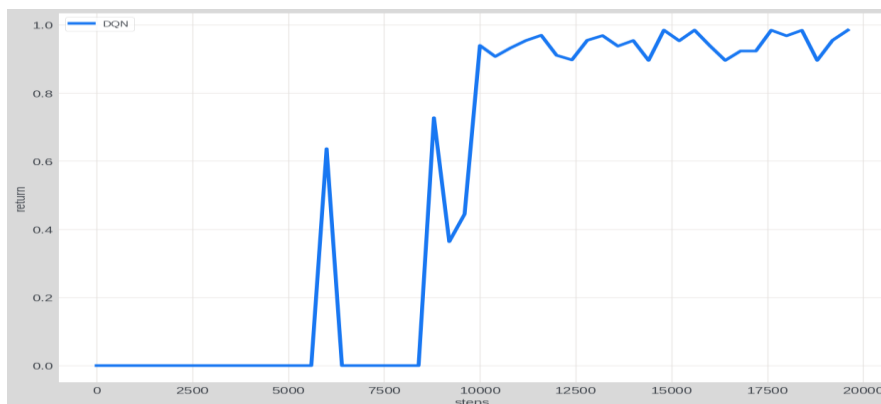
torch.save(info["return"], "DQN-return.pt")
plt.plot(record_period * np.arange(len(info["return"])), info["return"],
label="DQN")
plt.xlabel("steps")
plt.ylabel("return")
plt.legend()
plt.show()

```

---

<sup>11</sup> Agent

نمایش نتیجه :



با افزایش تعداد قدم ها، می توانیم افزایش در مقدار متوسط بازگشتی را مشاهده کنیم.