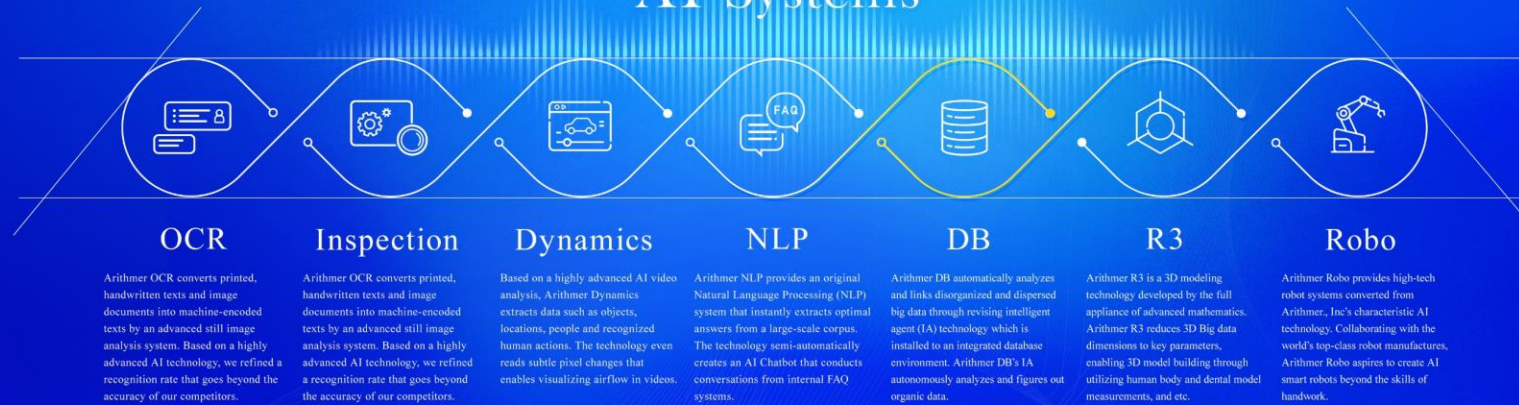# Recommendation Algorithm Using Reinforcement Learning

**Arithmer  DB Lu Juanjuan**

2020/09/15

- Lu Juanjuan
  - Graduated School
    - Tokyo Institute of Technology
      - Ishida Takashi Laboratory, Department of Computer Science , School of Computing
        Master research domain:
          Drug discovery by applying machine learning technologies

  - Current Job
    - Arithmer Inc.（Home page: https://arithmer.co.jp/en/）
      - Application of Machine Learning/ Data Analysis

# 1. Background

# 2. System Structure

# Background

Recommendation Algorithms:



COLLABORATIVE FILTERING ①

Read by both users

Similar users

Read by her, recommended to him!

Similar items

(item-based)

(user-based)

CONTENT-BASED FILTERING ②

Read by user

Similar articles

Recommended to user

[1]

Deep Learning Models

Predict: click or not

Model

Input data

[1]TONDJI, LIONEL NGOUPEYOU. "Web recommender system for job seeking and recruiting." (2018).

Artificial Intelligence

Machine Learning

"Machine" = Model

Neural network
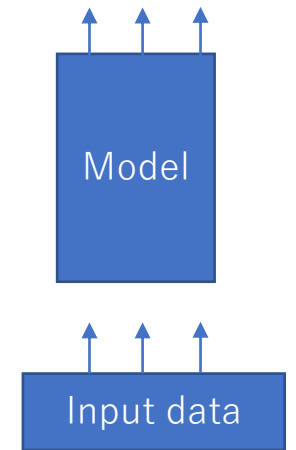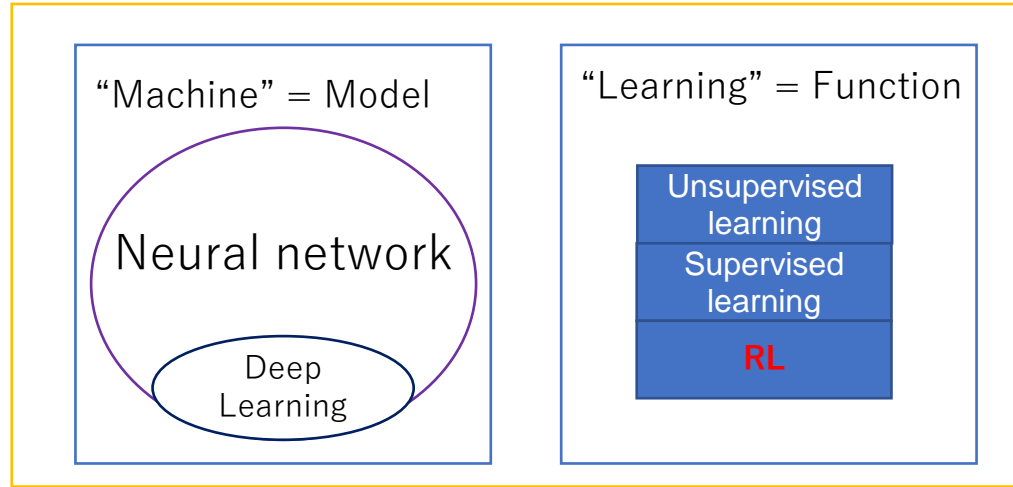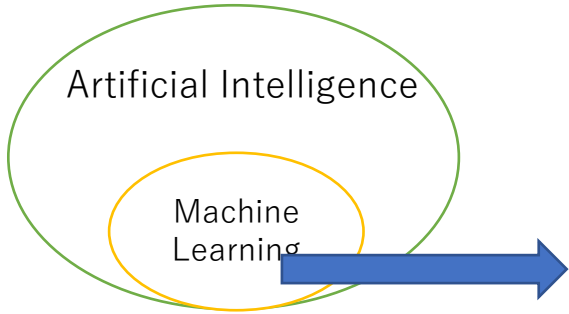
Deep Learning

"Learning" = Function

Unsupervised learning

Supervised learning

**RL**

[2]

S1: state,
a1,a2,a3,a4: actions

**Two major RL types:**

valued-based、policy-based

Q-learning: update **Q value** table

| action state | a1 | a2 | a3 | a4 |
|---|---|---|---|---|
| S1 | Q(S1, a1) | Q(S1, a2) | Q(S1, a3) | Q(S1, a4) |

$$Q(S,A) \leftarrow (1-\alpha)Q(S,A) + \alpha[R(S,a) + \gamma maxQ(S',a)]$$
$$a$$

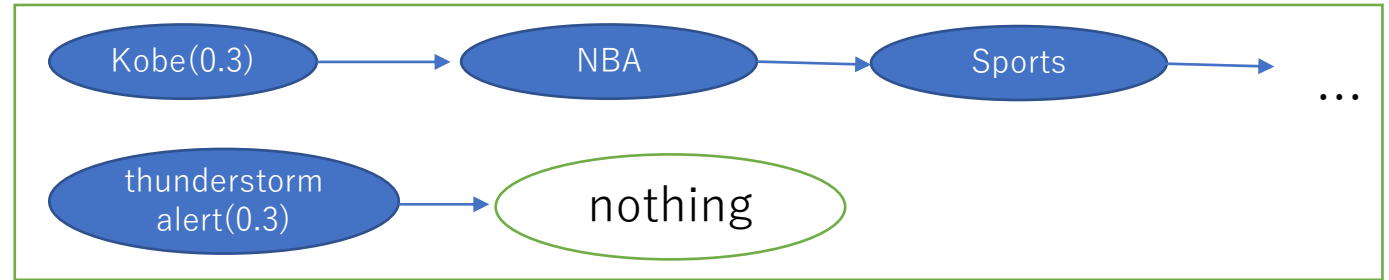| -1 | -1 | a1 -1 | -10 | -1 | -1 |
|---|---|---|---|---|---|
| -10 | a4 -1 | S-1 a2 -1 | -1 | -1 | -1 |
| -1 | -10 | a3 -1 | -10 | 20 | -1 |
| -1 | -10 | -1 | -10 | -10 | -1 |
| 0 | -1 | -1 | -1 | -1 | -1 |

Policy Gradient: update **policy** by gradient descent

$$\mathrm{E}_{\tau \sim \pi_\theta}[R(\tau)\nabla_\theta log\pi_\theta(\tau)]$$

[2]Kubo, Takahiro. *Paison De Manabu Kyoka Gakushu: Nyumon Kara Jissen Made*. Kodansha., 2019.

# Reasons:

1. **Long term rewards**

2. **Having some randomness**

Probability: [0.1, 0.2, 0.3, 0.4], not always the 4th item be chosen
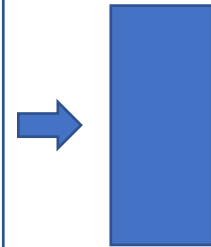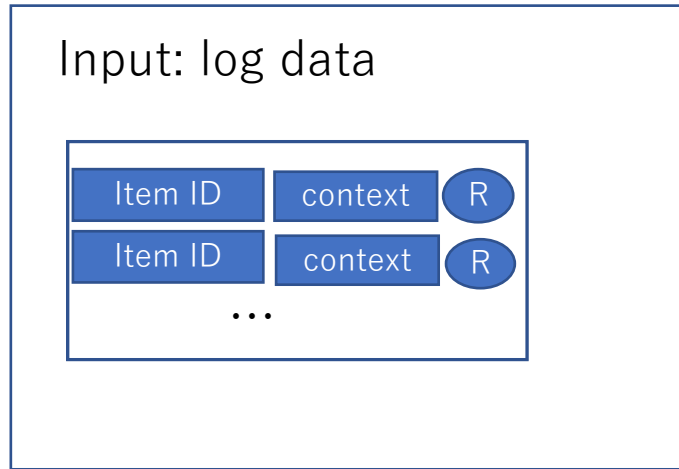


# Example:

1. **Policy Gradient based framework: being used to recommend videos.** [3]

   1. off-policy
   2. Continuous user state
   3. Experiment in live experiments

2. **DQN based framework: being used to recommend news.[4]**

3. **Critic-Actor based framework: being used to create a virtual environment like virtual Taobao.**

[3]Chen, Minmin, et al. "Top-k off-policy correction for a REINFORCE recommender system." *Proceedings of the Twelfth ACM International Conference on Web Search and Data Mining*. 2019.

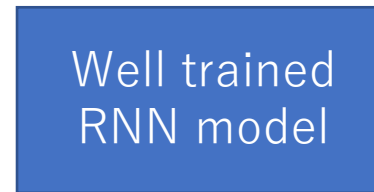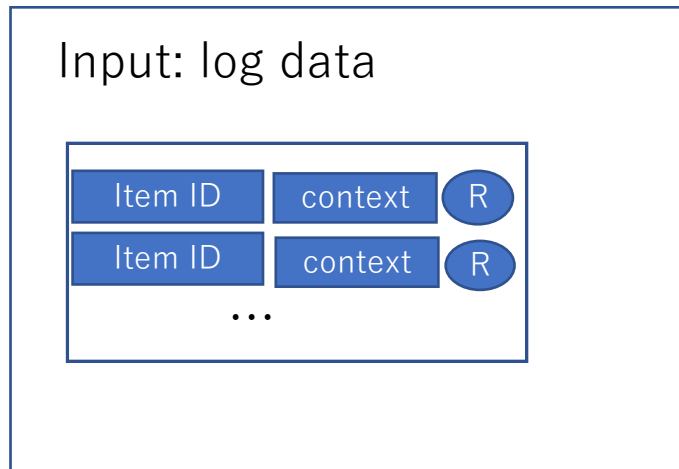[4]Zheng, Guanjie, et al. "DRN: A deep reinforcement learning framework for news recommendation." *Proceedings of the 2018 World Wide Web Conference*. 2018.

# System Structure
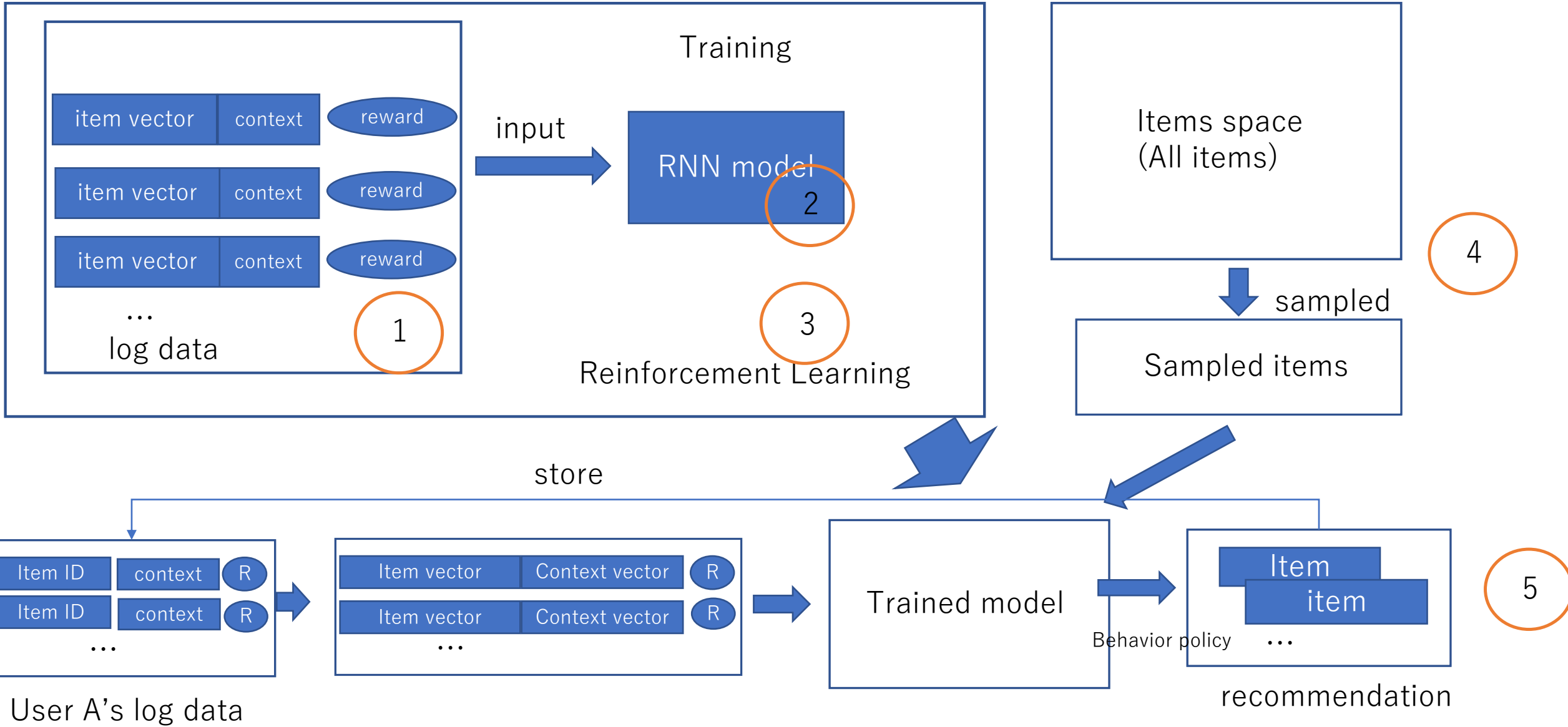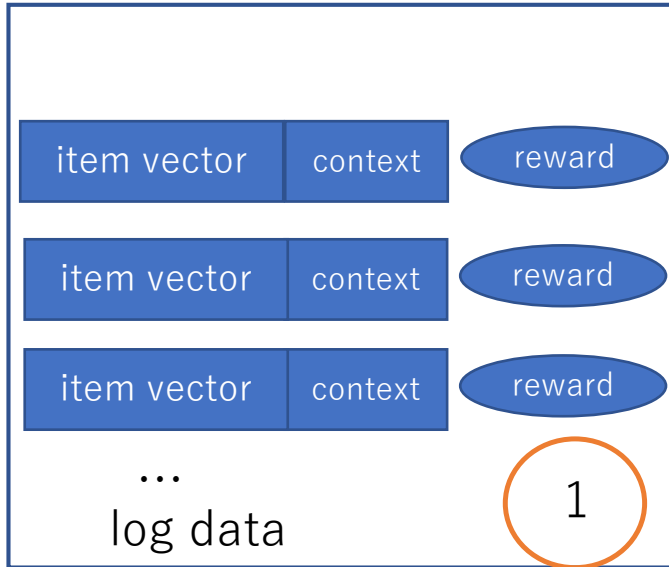
● Item vector:

Example：カジュアルコンフォート。【春夏生地】メリノウールにポリエステルを混紡した丈夫でしわになりにくい素材です。 48000。

Embedding: Word2vec/Bert

● Context data:

Example：timing、device

● Reward:
Example：1.click:  5 point,  2.buy:   15 point
3.non-feedback: 0 point

RNN model
2

s: state
A: whole item space
a: one item
$u_a$: item embedding + context vector
T: temperature(0~1)
$v_a$ : item embedding



Event t+1

Softmax $\pi_\theta$   Softmax $\beta_{\theta'}$

Block gradient

RELU

User state $s_{t+1}$   Label context

Recurrent Cell

Item embedding   Context

Event 1, 2,…, t

[3]

CFN cell

$$s_{t+1} = z_t \odot \tanh(s_t) + i_t \odot \tanh(W_a u_{a_t})$$
$$z_t = \sigma(U_z s_t + W_z u_{a_t} + b_z)$$
$$i_t = \sigma(U_i s_t + W_i u_{a_t} + b_i) \qquad [2]$$

$\pi_\theta$

$$\pi_\theta(a|s) = \frac{\exp(s^T v_a / T)}{\sum_{a' \in A} \exp(s^T v_{a'} / T)} \qquad [2]$$

$\beta_{\theta'} \ (behavior\ policy\ )$

$$\beta_{\theta'}(A|s) = \frac{\exp(s^T v_A / T)}{\sum_{a' \in A} \exp(s^T v_{a'} / T)}$$

RNN model

2

Event t+1

Softmax $\pi_\theta$

Softmax $\beta_{\theta'}$

Block gradient

RELU

User state $s_{t+1}$ | Label context

Recurrent Cell

Item embedding | Context

Event 1, 2,…, t

[3]

User State

S0 → CNF CELL → S1 → CNF CELL → S1 → … → St → CNF CELL → St+1

R0(!=0)          R1(==0)                    Rt(!=0)

Item embedding| context     Item embedding| context     …     Item embedding| context

a0                          a1                                at

**Ignoring non-reward item**

*S0 : [0,0,0,···,0]

RNN model
2

Event t+1

Softmax $\pi_\theta$    Softmax $\beta_{\theta'}$

**Block gradient**

RELU

User state $s_{t+1}$    Label context

Recurrent Cell

Item embedding    Context

Event 1, 2,…, t

[3]

$$\pi_\theta(at|st)$$

Softmax layer

Item embedding    User state

$$argmax(\beta_{\theta'}(A|s))$$

Softmax layer

Item embedding    User state

教師あり
でトレニ
ング

Reinforce algorithm:

Gradient

Policy

Trajectory: (s0,a0,s1,a1,..,sn,an)

$$\mathrm{E}_{\tau \sim \pi_\theta}[R(\tau)\nabla_\theta log\pi_\theta(\tau)]$$

Reward

Off policy:

**Important weight** of the off-policy-corrected gradient estimator

$$\sum_{\tau \sim \beta}[\sum_{t=0}^{|\tau|}\frac{\pi_\theta(a_t|s_t)}{\beta\,(a_t|s_t)}R_t\nabla_\theta log\pi_\theta(a_t|s_t)]$$

Top K:

$$\sum_{\tau \sim \beta} [\sum_{t=0}^{|\tau|} \frac{\alpha_\theta(a_t|s_t)}{\beta\,(a_t|s_t)} R_t \nabla_\theta log\alpha_\theta(a_t|s_t)]$$

$$\lambda_{K(s_t,\,a_t)} = \frac{\partial_\alpha(a_t|s_t)}{\partial_\pi(a_t|s_t)} = K(1 - \pi_\theta(a_t|s_t))^{K-1}$$

$$= \sum_{\tau \sim \beta} [\sum_{t=0}^{|\tau|} \frac{\pi_\theta(a_t|s_t)}{\beta\,(a_t|s_t)} \frac{\partial_\alpha(a_t|s_t)}{\partial_\pi(a_t|s_t)} R_t \nabla_\theta log\pi_\theta(a_t|s_t)]$$

## Final training expression:

$$\sum_{\tau \sim \beta} [\sum_{t=0}^{|\tau|} \frac{\pi_\theta(a_t|s_t)}{\beta\,(a_t|s_t)} K(1 - \pi_\theta(a_t|s_t))^{K-1} R_t \nabla_\theta log\pi_\theta(a_t|s_t)]$$

Items space
(All items)

**During server time:**

sampled

④

Efficient approximate nearest neighbor-based systems

Sampled items

Web page

| item1 ♥ | item2 | item3 | item4 | item5 ♥ |
| item6 | item7 | Item8 ♥ | item9 | item10 |
| item11 ♥ | item12 | item13 | item14 | item15 |
...

*30 popular items from each category

Step 3

*Event t+1*

Softmax $\pi_\theta$   Softmax $\beta_{\theta'}$

Block gradient

RELU

Step 1

User state $s_{t+1}$   Label context

Recurrent Cell

Item embedding   Context

*Event 1, 2,…, t*

[3]

Step 2

Items space
(All items)   sampled

Sampled items

Step1: Choosing 10 items and then get user's state vector.

Step2: Sampling items from items space.

Step3: Calculating recommendation probability of all sampled items.

Step4: Randomly recommend K items with recommendation probability.

Step5: Storing recommended item info , context info and users' feedback.

Step 3

Step 1

Step 2

Log data

[3]

Items space (All items)

sampled

Sampled items

Step1: Getting user's state vector by inputting log data.

Step2: Sampling items from items space.

Step3: Calculating recommendation probability of all sampled items.

Step4: Randomly recommend K items with recommendation probability.

Step5: Storing recommended item info , context info and users' feedback.

**Arithmer 株式会社**

〒106-6040
東京都港区六本木一丁目 6 番 1 号 泉ガーデンタワー 38/40F ( 受付 )
03-5579-6683
https://arithmer.co.jp/