

Georgios Frangias

August 2023

# Algorithm 1 On-policy SLATEQ for Live Experiments

## 1: **Parameters:**

- $T$  : the number of iterations.
- $M$  : the interval to update label network.
- $\gamma$  : discount rate.
- $\theta_{\text{main}}$  : the parameter for the main neural network.
- $\bar{Q}_{\text{main}}$  : that predicts items' long-term value.
- $\theta_{\text{label}}$  : the parameter for the label neural network  $\bar{Q}_{\text{label}}$ .
- $\theta_{\text{pctr}}$  : the parameter for the neural network that predicts items' pCTR.

## 2: **Input:**

- $D_{\text{training}} = (s, A, C, L_{\text{myopic}}, s', A')$  : the training data set.
- $s$  : current state features
- $A = (a_1, \dots, a_k)$  : recommended slate of items in current state;  $a_i$  denotes item features
- $C = (c_1, \dots, c_k)$  :  $c_i$  denotes whether item  $a_i$  is clicked
- $L_{\text{myopic}} = (l_{\text{myopic}}^1, \dots, l_{\text{myopic}}^k)$  : myopic (immediate) labels
- $s'$  : next state features
- $A' = (a'_1, \dots, a'_k)$  : recommended slate of items in next state.

3: **Output:** Trained Q-network  $\bar{Q}_{\text{main}}$  that predicts items' long-term value.

4: **Initialization:**  $\theta_{\text{label}} = 0, \theta_{\text{main}}$  randomly,  $\theta_{\text{pctr}}$  randomly

5: **for**  $i = 1 \dots T$  **do**

6:     **if**  $i \bmod M = 0$  **then**

7:          $\theta_{\text{label}} \leftarrow \theta_{\text{main}}$

8:     **for** each example  $(s, A, C, L_{\text{myopic}}, s', A') \in D_{\text{training}}$  **do**

9:         **for** each item  $a_i \in A$  **do**

10:             update  $\theta_{\text{pctr}}$  using click label  $c_i$

11:             **if**  $a_i$  is clicked **then**

12:                 probability:  $p \text{CTR}(s', a'_i, A') \leftarrow p \text{CTR}(s', a'_i) / \sum_{a'_i \in A} p \text{CTR}(s', a'_i)$

13:                 LTV label:  $l_{\text{ltv}}^i \leftarrow l_{\text{myopic}}^i + \sum_{a'_i \in A'} p \text{CTR}(s', a'_i, A') \bar{Q}_{\text{label}}(s', a'_i)$

14:             update  $\theta_{\text{main}}$  using LTV label  $l_{\text{ltv}}^i$

Taken from arXiv:1905.12767