# Towards More Realistic User Long Term Engagement Modeling in Recommender Systems

## Xu Chen

# CONTENTS

1. Interactive Nature

2. Long-term Utilities

**State (S)**

**Reward (R)**

Environment

Agent

**Action (A)**

**Reward** (User's feedback for the item)

User

Recommender System

**Action** (Recommending an item)

Recommender -> Agent

User -> Environment

Recommend an item -> Action
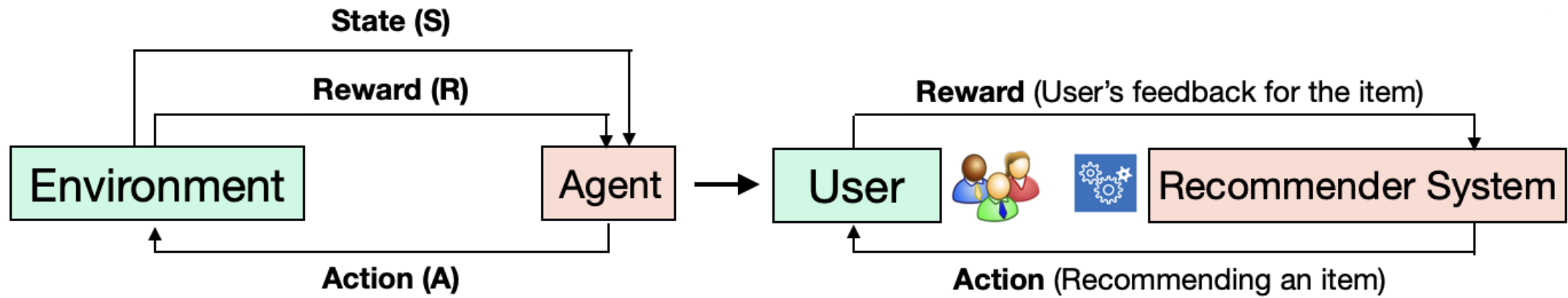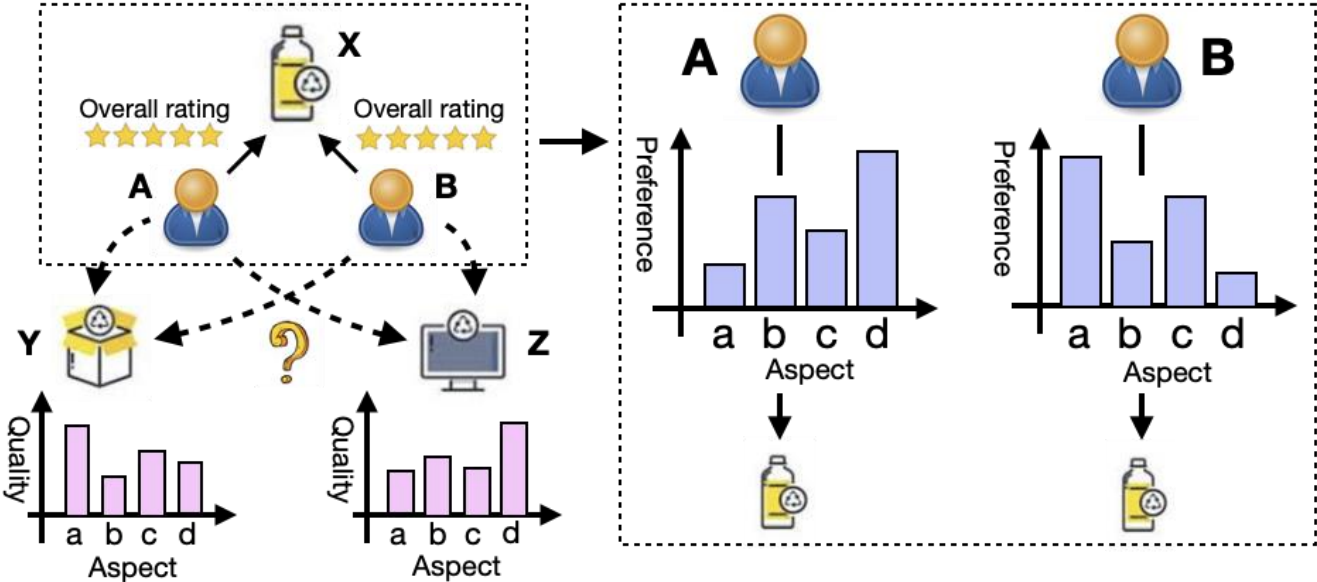
User's historical behaviour -> State

User's rating on the item -> Reward

How to consider the characters of the user preference? E.g., Diversity, Dynamic, …

Diverse user preference

Different aspect preferences are not always aligned!

**Definition 1. Pareto dominance.** Suppose we have two parameters $\theta_A$ and $\theta_B$, we say $\theta_A$ can dominant $\theta_B$ (denoted by $\theta_A > \theta_B$), if and only if $\mathcal{L}_i(\theta_A) \leq \mathcal{L}_i(\theta_B)$, $\forall i \in \{1, 2, ...M\}$ and $\mathcal{L}_i(\theta_A) < \mathcal{L}_i(\theta_B)$, $\exists i \in \{1, 2, ...M\}$.

**Definition 2. Pareto efficiency.** For a parameter $\theta^*$, if there is no other $\hat{\theta}$, such that $\hat{\theta} > \theta^*$, then we say $\theta^*$ is a Pareto efficient solution.

Critic learning:

$$\arg\min_{\boldsymbol{\phi}_m} \sum_{i=1}^{N} (y_{i,m} - Q_m(s_i, a_i|\boldsymbol{\phi}_m))^2, \ m = 1, 2, ...M$$

Actor learning:

$$l(\boldsymbol{\theta}) = -\sum_{m=1}^{M} w_m \sum_{i=1}^{N} Q_m(s_i, \mu(s_i|\boldsymbol{\theta}))$$

$$L(\boldsymbol{\theta}) = \sum_{i=1}^{N} y_i \log \sigma(\boldsymbol{q}_o^T \mu(s_i|\boldsymbol{\theta})) + (1 - y_i) \log(1 - \sigma(\boldsymbol{q}_o^T \mu(s_i|\boldsymbol{\theta}))) \longrightarrow L(\boldsymbol{\theta}) = \tilde{Q}((o, y_i), \mu(s_i|\boldsymbol{\theta}))$$

$$\min_{\boldsymbol{w}} || \sum_{m=1}^{M} w_m \nabla_{\boldsymbol{\theta}} \sum_{i=1}^{N} Q_m(s_i, \mu(s_i|\boldsymbol{\theta})) ||_2^2$$
$$s.t. \ \boldsymbol{e}_k^T \boldsymbol{w} \geq b_k, \ \forall \ k \in [1, K]$$
$$\boldsymbol{1}^T \boldsymbol{w} = 1, \ w_m \geq 0, \ \forall \ m \in [1, M]$$

**Theorem 1.** *If $\boldsymbol{w}$ is determined by solving the quadratic programming (QP) problem of (5), then either one of the following holds:*
*i) The solution to the optimization problem is 0, then the local Pareto efficient solution is achieved.*
*ii) $\boldsymbol{d} = \sum_{m=1}^{M} w_m \nabla_{\boldsymbol{\theta}} \sum_{i=1}^{N} Q_m(s_i, \mu(s_i|\boldsymbol{\theta}))$ is a gradient direction which does not decrease any Q function.*

**Algorithm 1:** Pareto Deterministic Policy Gradient

1 Initialize Actor parameter $\theta$ and Target Actor parameter $\theta' \leftarrow \theta$.

2 Initialize Critic parameter $\phi_m$ and Target Critic parameter $\phi'_m \leftarrow \phi_m, \forall m \in [1, M]$.

3 Initialize Pareto weights $w = \{\frac{1}{M+1}, \frac{1}{M+1}, ..., \frac{1}{M+1}\}$ and replay buffer $B$.

4 **for** *episode number in [1, K]* **do**

5    **i) Trajectory Generation**

6    Get start state $s_1$

7    **for** *step t in [1, T]* **do**

8      Select an action according to $a_t = \mu(s_t|\theta) + N_t$, $N_t$ is an exploration noise.

9      Execute $a_t$ to obtain the new state $s_{t+1}$ and the reward vector $r_t = \{r_{t,1}, r_{t,2}, ..., r_{t,M}\}$.

10      Push $\{s_t, a_t, r_t, s_{t+1}\}$ into the replay buffer $B$

11    **end**

12    **ii) Update Critic**

13    Sample Z instances $\{s_i, a_i, r_i, s_{i+1}\}$ from $B$

14    **for** *critic m in [1, M]* **do**

15      **for** *i in [1, Z]* **do**

16        Compute $y_i = r_{i,m} + \gamma Q_m(s_{i+1}, \mu(s_{i+1}|\theta')|\phi'_m)$.

17      **end**

18      $\phi_m \leftarrow \phi_m - \alpha_\phi \nabla_{\phi_m}\{\frac{1}{Z}\sum_{i=1}^{Z}(y_i - Q_m(s_i, a_i|\phi_m))^2\}$.

19    **end**

20    **iii) Update Pareto Weight**

21    **for** *i in [1, Z]* **do**

22      **for** *m in [1, M]* **do**

23        $p_{i,m} = \nabla_a Q_m(s, a)|_{s=s_i, a=\mu_\theta(s_i)} \nabla_\theta \mu(s|\theta)|_{s=s_i}$.

24      **end**

25      $\tilde{p}_i = \nabla_a \tilde{Q}(s, a))|_{s=(o, y_i), a=\mu_\theta(s_i)} \nabla_\theta \mu(s|\theta)|_{s=s_i}$.

26    **end**

27    Update $w = \{w_1, w_2, ..., w_M, \tilde{w}\}$ by Solving (5).

28    **iv) Update Actor**

29    $d = \frac{1}{Z}\sum_{i=1}^{Z}\sum_{m=1}^{M} w_m p_{i,m} + \tilde{w}\tilde{p}_i$.

30    $\theta \leftarrow \theta + \alpha_\theta d$.

31    $\theta' \leftarrow \tau\theta + (1 - \tau)\theta'$.

32    $\phi'_m \leftarrow \tau\phi_m + (1 - \tau)\phi'_m \quad \forall m \in [1, M]$.

33 **end**

$$G = \left\|\mathbb{E}_{B_i}\left[\sum_{m=1}^{M+1} w_m(B_i)(\frac{1}{Z}\sum_{s_b \in B_i} f_m(s_b; \theta) - \mathbb{E}_s[f_m(s; \theta)])\right]\right\|$$

**Theorem 2.** *Suppose i) $\nabla_a Q_m(s, a)$ and $\nabla_\theta \mu(s|\theta)$ are bounded by $X_m$ and $Y$, that is, $\|\nabla_a Q_m(s, a)\|_2 \leq X_m$ and $\|\nabla_\theta \mu(s|\theta)\|_2 \leq Y$. ii) The batched gradient of the action-value function for each objective is unbiased, that is: $\mathbb{E}_{B_i}[\frac{1}{Z}\sum_{s_b \in B_i} f_m(s_b; \theta)] = \mathbb{E}_s[f_m(s; \theta)]$. iii) $f_m(s_b; \theta)$ follows a normal distribution $\mathcal{N}(\mathbb{E}_s[f_m(s; \theta)], \sigma^2 I)$, where $I \in \mathbb{R}^{d\times d}$ is an identity matrix and $\sigma$ is a scalar. Then we have:*

$$G \leq \sum_{m=1}^{M+1} \mathbb{E}_{B_i}\left[w_m(B_i)(|\frac{1}{Z}\sum_{s_b \in B_i} f_m(s_b; \theta) - \mathbb{E}_s[f_m(s; \theta)]|)\right] \quad (14)$$

$$\leq \frac{X_{m*} Y \sqrt{d}}{\sqrt{Z}} \quad (15)$$

*where $m^* = \arg\max_m |\frac{1}{Z}\sum_{s_b \in B_i} f_m(s_b; \theta) - \mathbb{E}_s[f_m(s; \theta)]|$.*
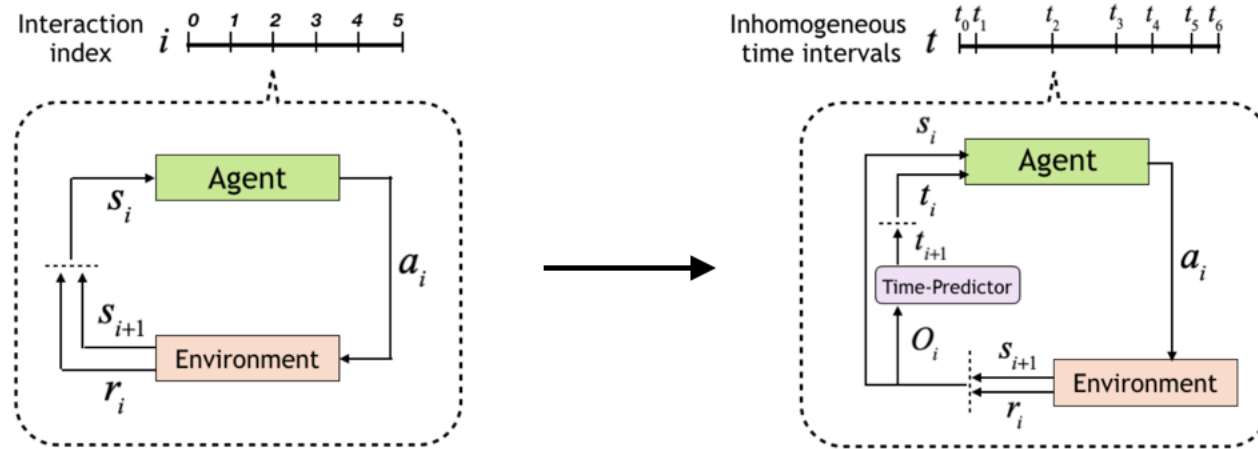
**Weight-reuse mechanism.** In this method, we introduce a container $\boldsymbol{W} \in \mathbb{R}^{L \times (M+1)}$ for storing previously derived Pareto weights. For each training batch $\boldsymbol{B}_i$, $\boldsymbol{w} \in \mathbb{R}^{M+1}$ is not always computed by solving problem (5). We firstly check the weights in the container:

(1) If there is a candidate $\boldsymbol{w}^* \in \boldsymbol{W}$, such that its corresponding $\boldsymbol{d}^* = \sum_{m=1}^{M+1} w_m^* \nabla_{\boldsymbol{\theta}} \left( \frac{1}{Z} \sum_{s_b \in \boldsymbol{B}_i} Q_m(s_b, \mu(s_b|\boldsymbol{\theta})) \right)$ can increase all the Q functions, that is, $(\boldsymbol{d}^*)^T \nabla_{\boldsymbol{\theta}} \left( \frac{1}{Z} \sum_{s_b \in \boldsymbol{B}_i} Q_m(s_b, \mu(s_b|\boldsymbol{\theta})) \right) > 0, \forall m \in [1, M+1]$, then we set $\boldsymbol{w} = \boldsymbol{w}^{*3}$. Since the weights in $\boldsymbol{W}$ is not derived from $\boldsymbol{B}_i$, the bias $G$ becomes 0 at this moment.

(2) If there is no such weight in $\boldsymbol{W}$, we solve problem (5) to derive $\boldsymbol{w}$, which is then pushed into the container for future "reuse". In this scenario, $G$ is not 0, which is bounded by equation (15).

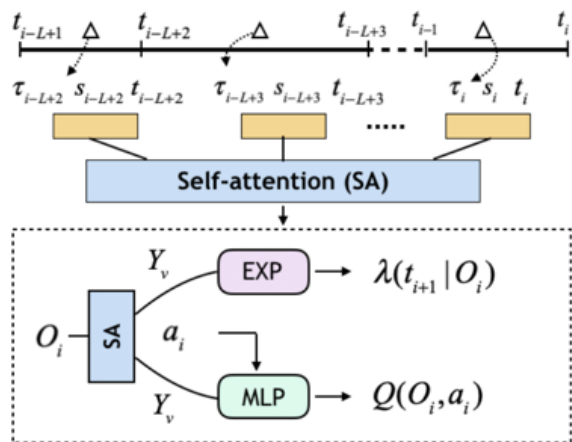# Modelling User Dynamic Preference

Dynamic user preference

Inhomogeneous DQN

Trajectory

$$\tau = \{s_0, a_0, t_0, s_1, a_1, t_1, \ldots, s_T, a_T, t_T\}$$

Objective

$$J = \mathbb{E}_{\tau \sim p_{\theta, \psi}(\tau)}\left[\sum_{i=0}^{T} \gamma^i \kappa(t_0 - t_i) r(s_i, a_i, t_i))\right]$$

Bellman Operator

$$(TQ)(s_i, a_i, t_i)$$
$$= r(s_i, a_i, t_i) + \sum_{s_{i+1} \in \mathcal{S}} \mathcal{P}(s_{i+1}|s_i, a_i)\{\mathbb{E}_{t_{i+1} \sim \mathcal{T}(\cdot|O_i)}[\gamma \kappa(t_i - t_{i+1})$$
$$\max_{a_{i+1} \in \mathcal{A}} Q(s_{i+1}, a_{i+1}, t_{i+1})]\}$$

**Theorem 1** ($T$ is a contractive operator.). *Let $Q_1$ and $Q_2$ be two value functions. Then, the Lipschitz condition $\|TQ_1 - TQ_2\|_\infty \leq \alpha\|Q_1 - Q_2\|_\infty$ holds, where $\alpha \in [0, 1)$ is a constant.*

$$Q(O_i, a_i | \boldsymbol{\phi}_Q, \boldsymbol{\phi}_Y) = f_Q^n(f_Q^{n-1}(\dots(f_Q^1(Y_v \boldsymbol{w}_3, \boldsymbol{p}_i))\dots))$$

$$\lambda^*(t_{i+1} | \boldsymbol{\phi}_\lambda, \boldsymbol{\phi}_Y) = \lambda(t_{i+1} | O_i; \boldsymbol{\phi}_\lambda, \boldsymbol{\phi}_Y)$$

$$= \exp\big(\underbrace{\boldsymbol{w}_1^T Y_v \boldsymbol{w}_2}_{A} + \underbrace{\boldsymbol{w}_t(t_{i+1} - t_i)}_{B} + \underbrace{b}_{C}\big)$$

$$\boldsymbol{x}_j = [\boldsymbol{e}_j; \boldsymbol{\tau}_j] + \boldsymbol{v}_j$$

$$\boldsymbol{Q} = \boldsymbol{X}^T \boldsymbol{W}_Q, \boldsymbol{K} = \boldsymbol{X}^T \boldsymbol{W}_K, \boldsymbol{V} = \boldsymbol{X}^T \boldsymbol{W}_V,$$

$$\boldsymbol{X}_1 = \text{SOFTMAX}(\frac{\boldsymbol{Q}\boldsymbol{K}^T}{\sqrt{d_K}})\boldsymbol{V},$$

$$\boldsymbol{Y}_1 = \boldsymbol{W}_2^F \text{ReLU}(\boldsymbol{W}_1^F \boldsymbol{X}_1^T + \boldsymbol{b}_1^F) + \boldsymbol{b}_2^F,$$

$$L(\phi_Q, \phi_Y, \phi_\lambda) = E_{(O_i, a_i, r_i, s_{i+1}, t_{i+1})}[(y_i - Q(O_i, a_i))^2]$$

$$= \int_{O_i} \underbrace{p(O_i)}_{A} \underbrace{E_{(a_i, r_i, s_{i+1}, t_{i+1}|O_i)}[(y_i - Q(O_i, a_i))^2]}_{B} \, dO_i$$

A

$$L_P(\phi_\lambda, \phi_Y) = \sum_{O_i} \sum_{j=1}^{i} \log p(t_j|O_{j-1}) = \sum_{O_i} \sum_{j=1}^{i} \log f^*(t_j|O_{j-1})$$

$$= \sum_{O_i} \sum_{j=1}^{i} \{\log \lambda^*(t_j|\phi_\lambda, \phi_Y) - \int_{t_{j-1}}^{t_j} \lambda^*(\tau|\phi_\lambda, \phi_Y)d\tau\}$$

B

$$L_Q(\phi_Q, \phi_Y) = \sum_{D} (y_i - Q(O_i, a_i|\phi_Q, \phi_Y))^2$$

**Theorem 2** (Finite-time bound for IDQN). *Suppose we have $K$ training iterations in the optimization process, and for the $k$th iteration, the Q-network is updated from $Q_{k-1}$ to $Q_k$ as follows:*

$$Q_k = \arg \min_{f \in \mathcal{F}} \sum_{i=1}^{N} [y_i - f(O_i, a_i)]^2$$
$$y_i = r_i + \gamma \kappa(t_i - t_{i+1}) \max_{a_{i+1}} Q_{k-1}(O_{i+1}, a_{i+1})$$

(14)

*where we have $N$ training samples, and $(O_i, a_i)$ in each instance is drawn from a distribution $\sigma$.*

*Suppose (i) $k(m; \mu, \sigma)$ is the concentration coefficient as defined in [37], where $\mu$ is the distribution on $(O_i, a_i)$ after $m$ MDP steps, and $k(m; \mu, \sigma)$ measures the similarity between $\mu$ and $\sigma$. (ii) $l^k$ is the minimum time interval spanning $k$ steps in all trajectories, that is, $l^k = \min_{i,\tau}(t_{i+k}^{\tau} - t_i^{\tau})$, where $t_i^{\tau}$ is the agent-environment interaction time for the $i$th step in trajectory $\tau$ (iii) $e_{k+1} = TQ_k - Q_{k+1}$ and $s = \max_i \|e_i\|_{\sigma}$. We assume (i) the immediate reward is bounded by $R_{max}$, (ii) $\sum_{m \geq 1}[\gamma^{m-1} m k(m; \mu, \sigma)]^2 \leq \phi_{\mu,\sigma}$, and (iii) $\sum_{k=0}^{\infty} \kappa(-l^k)^2 \leq \phi_{\kappa}$.*

*Let $\pi_K$ be the one-step greedy policy of $Q_K$, and $Q^{\pi_K}$ be the action-value function corresponding to $\pi_K$. $Q^*(O, a) = \sup_{\pi} Q^{\pi}(O, a)$ is the optimal Q-function. Then, the upper bound of the error between $Q^*$ and $Q^{\pi_K}$ is:*

$$|Q^* - Q^{\pi_K}|_{1,\mu} \leq 2\gamma s(\phi_{\mu,\sigma} \phi_{\kappa})^{\frac{1}{2}} + \frac{4\gamma^{K+1} R_{max}}{(1-\gamma)^2}$$

(15)

# Outlook

What if the logged user preference is biased?

    Counterfactual trajectory generation

    Debiased user preference learning

    ...

How to build high reliable simulator?

    Not limited to RL-based Recsys

    Potential to promote the Recsys research

    ...

# Thanks & QA