

# Self-Modulating Attention in Continuous Time Space with Applications to Sequential Recommendation

Chao Chen, Haoyu Geng,  
Nianzu Yang, Junchi Yan,  
Daiyue Xue, Jianping Yu, Xiaokang Yang



## Abstract

User interests are usually dynamic in the real world, which poses both theoretical and practical challenges for learning accurate preferences from rich behavior data. Among existing user behavior modeling solutions, attention networks are widely adopted for its effectiveness and relative simplicity. Despite being extensively studied, existing attentions still suffer from two limitations: i) conventional attentions mainly take into account the spatial correlation between user behaviors, regardless the distance between those behaviors in the continuous time space; and ii) these attentions mostly provide a dense and undistinguished distribution over all past behaviors then attentively encode them into the output latent representations. This is however not suitable in practical scenarios where a user's future actions are relevant to a small subset of her/his historical behaviors. In this paper, we propose a novel attention network, named self-modulating attention, that models the complex and non-linearly evolving dynamic user preferences. We empirically demonstrate the effectiveness of our method on top-N sequential recommendation tasks, and the results on three large-scale real-world datasets show that our model can achieve state-of-the-art performance.

## Self-modulating Attention

The motivation for designing attention in continuous time space is to directly inform the correlation  $p(v_{ij} | H_t)$  with the expected number of occurrence  $E[N(t, t+dt) | H_{t-}, v_{ij}]$ , conditional on the history  $H_{t-}$ , where  $H_{t-} = H_t \cup \{t_j + 1/ \in (t_j, t)\}$  and  $N(t, t+dt) \in \{0, 1\}$  denotes the number of occurrences for item  $i_j$  in an infinitesimal interval.

$$\begin{aligned} & \mathbb{E}_{p(\mathbf{v}|\mathcal{H}_t)} [\mathbb{E}[N(t, t+dt) | \mathcal{H}_{t-}, \mathbf{v}] \mathbf{v}] \\ & \stackrel{(a)}{=} \mathbb{E}_{p(\mathbf{v}|\mathcal{H}_t)} [p(t_{j+1} \in [t, t+dt) | \mathcal{H}_{t-}, \mathbf{v}) \mathbf{v}] \\ & \stackrel{(b)}{=} \mathbb{E}_{p(\mathbf{v}|\mathcal{H}_t)} \frac{p(t_{j+1} \in [t, t+dt) | \mathcal{H}_t, \mathbf{v})}{p(t_{j+1} \notin (t_j, t) | \mathcal{H}_t, \mathbf{v})} \mathbf{v} \\ & \stackrel{(c)}{=} \mathbb{E}_{p(\mathbf{v}|\mathcal{H}_t)} \frac{f(t | \mathcal{H}_t, \mathbf{v}) dt}{S(t | \mathcal{H}_t, \mathbf{v})} \mathbf{v} \\ & \stackrel{(d)}{=} \mathbb{E}_{p(\mathbf{v}|\mathcal{H}_t)} \mathbf{v} \lambda^*(t | \mathcal{H}_t, \mathbf{v}) dt \end{aligned}$$

The dynamic processes will be characterized by conditional intensity function  $\lambda^*(t | H_t)$ . The generalized self-modulating attention could be formulated as:

$$\widehat{\text{Att}}(\mathbf{q}, | \mathcal{H}_t) = \mathbb{E}_{p(\mathbf{v}|\mathcal{H}_t)} \mathbf{v} \lambda^*(t | \mathcal{H}_t, \mathbf{v})$$

## Self-modulating Layer (SMLayer)

We introduce a specific implementation, self-modulating layer (SMLayer) with application to sequential recommendation.

The widely used self-attention in Sequential recommendation has following form:

$$\begin{aligned} \mathbf{X} &= \text{concat}([\mathbf{Y}, \mathbf{Z}]) \\ \mathbf{Q} &= \mathbf{XW}^Q, \quad \mathbf{K} = \mathbf{XW}^K, \quad \mathbf{V}^{\text{seq}} = \mathbf{XW}^V \\ \mathbf{H} &= \text{softmax}\left(\frac{\mathbf{QK}^T}{\sqrt{d}}\right) \mathbf{V}^{\text{seq}}. \end{aligned}$$

where all  $W$  are projection matrices,  $Q, K, V_{\text{seq}}$  are separately the query, key and values matrices obtained by different transformations of the input  $X$ , and  $H$  is the output representations of conventional attentions. Further we present the formulation of conditional intensity function:

$$\begin{aligned} \mathbf{g}_k(t) &= \sigma\left(\underbrace{\mathbf{W}_k^G \mathbf{h}(t_j)}_{\text{Endogenous}} + \underbrace{\mathbf{b}_k^G(t - t_j)}_{\text{Exogenous}}\right) \\ \lambda(t | \mathcal{H}_t, \mathbf{v}_k) &= f_k(\mathbf{w}_k^T \mathbf{g}_k(t) + \mu_k) \end{aligned}$$

where  $g_k(t)$  signifies continuous time temporal dynamics, and  $f_k$  is the softplus activation function:

$$f_k(x) = \phi_k \log(1 + \exp(x/\phi_k))$$

## Continuous Time Regularization (CTReg)

In preference learning, the model parameters are usually optimized by minimizing the reconstruction error. The only source of supervised signal is from the behavior data (i.e.,  $R$ ) that is independent of time. Intensity function learned in protocol above might probably diverge from the complex continuous-time patterns contained in the data. CTReg (continuous time regularization) is defined by:

$$R(\Theta) = \sum_{j=1}^L \log \lambda_k(t_j | \mathcal{H}_j) - \int_{t_1}^{t_L} \lambda(t_j | \mathcal{H}_j) dt$$

The final objective function is minimization of empirical risk with CTReg as:

$$\min_{\Theta} \ell(\mathbf{R}, \hat{\mathbf{R}}) - \gamma \mathbb{E}_{u \in [1, m]} R(\Theta; u)$$

## Theoretical Results

Generalization Bound. Suppose that the loss function is  $L$ -Lipschitz, and for the estimate  $R$  on an random example set  $\Omega$ , we bound  $\rho|\Omega| = \sup_{P \in \mathcal{F}} \sum_{(u,i) \in \Omega} \|P_{i,\cdot} - \cdot\|_0$  and  $\mu = \sup_{(u,i,k) \in \Omega} |P_{u,k}(VB)_{k,j}|$ , then with probability at least  $1 - \delta$ , we have the bound:

$$\mathbb{E}[\ell(\mathbf{R}, \hat{\mathbf{R}})] \leq \mathbb{E}_{\Omega}[\ell(\mathbf{R}, \hat{\mathbf{R}})] + \mathcal{O}\left(L\mu\sqrt{\frac{C_p \ln|\Omega|}{|\Omega|}} + \sqrt{\frac{\ln(1/\delta)}{|\Omega|}}\right)$$

where  $C = d(m+n) \log(48emn)$

## Empirical Results

### Quantitative Analysis

We present experimental results on three large-scale datasets (Amazon, Koubei, Tmall) against sequential SOTA recommendation baselines.

Model	Amazon		Koubei		Tmall	
	HR	NDCG	HR	NDCG	HR	NDCG
SHAN	0.19250	0.11724	0.28150	0.20256	0.37316	0.25840
DIN	0.21955	0.13443	0.32665	0.24186	0.48460	0.33855
GRU4REC	0.24380	0.15822	0.32655	0.27052	0.46877	0.33746
SASREC	0.25595	0.16131	0.35455	0.27070	0.50433	0.34326
SASREC+	0.25820	0.16204	0.35690	0.27148	0.50607	0.34328
SASREC w/ ours.	<b>0.26545</b>	<b>0.16529</b>	<b>0.36235</b>	<b>0.28083</b>	<b>0.51347</b>	<b>0.35811</b>

Comparison to state-of-the-art Baselines on three benchmark Datasets

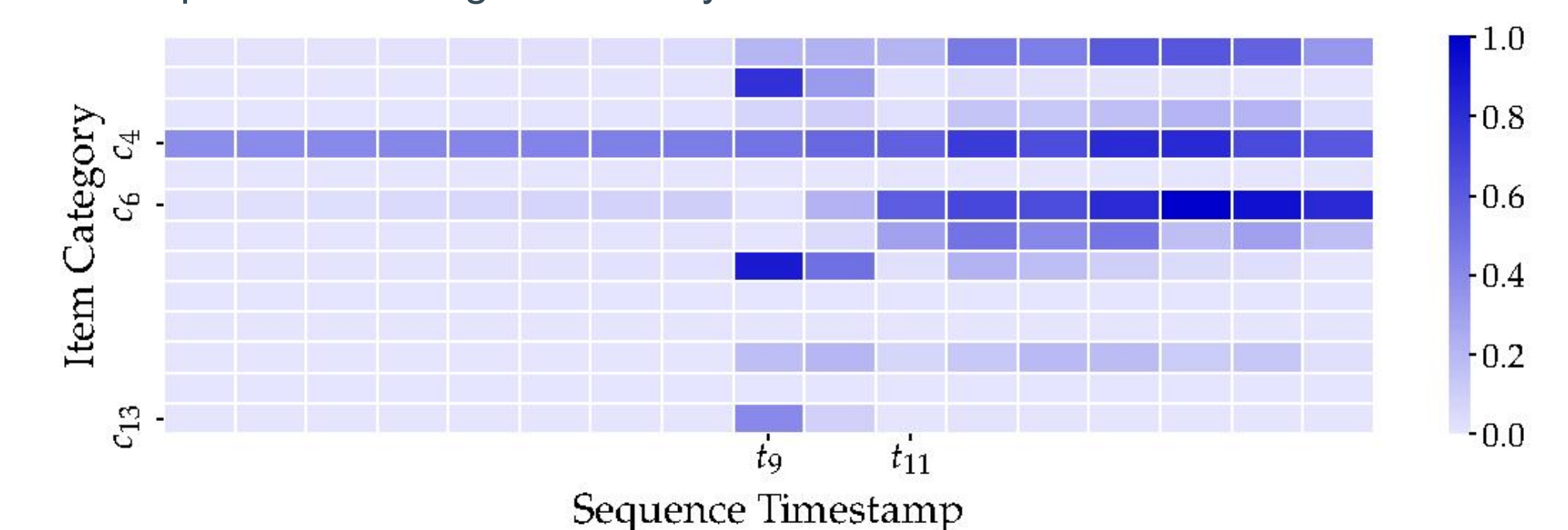
We further conduct ablation study to equip generalized self-modulating attention on DIN (Deep interest network) with proposed SMLayer and CTReg (Continuous time regularization) to validate our methods:

Dataset	DIN			SASREC		
	Origin	+SMLayer	+CTReg	Origin	+SMLayer	+CTReg
HR	Amazon	0.21955	<b>0.22065</b>	0.21985	0.25595	<b>0.26545</b>
	Koubei	0.32665	<b>0.33940</b>	0.33780	0.35455	<b>0.36194</b>
	Tmall	0.48460	<b>0.49033</b>	<b>0.49157</b>	0.50433	<b>0.51218</b>
NDCG	Amazon	0.13443	<b>0.13383</b>	0.13296	0.16131	0.16475
	Koubei	0.24186	<b>0.25444</b>	0.25411	0.27070	0.27862
	Tmall	0.33855	0.34580	<b>0.35062</b>	0.34326	0.35214

Ablation Study

### Qualitative Analysis

Example of user preference intensities on Koubei. The user usually purchases the items from category  $c_4$  for all the time, and he/she starts to repeatedly buy the items from category  $c_6$  after time  $t_9$ . The darker color corresponds to the higher intensity



### References

- Kang, W.-C. and McAuley, J. Self-attentive sequential recommendation. In Proceedings of the IEEE International Conference on Data Mining (ICDM '18), pp. 197–206. IEEE, 2018
- Xu, D., Ruan, C., Kumar, S., Korpeoglu, E., and Achan, K. Self-attention with functional time representation learning. In Proceedings of the Annual Conference on Neural Information Processing Systems (NeurIPS '19), pp. 15915–15925, 2019
- Zhou, G., Zhu, X., Song, C., Fan, Y., Zhu, H., Ma, X., Yan, Y., Jin, J., Li, H., and Gai, K. Deep interest network for click-through rate prediction. In Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD '18), pp. 1059–1068, 2018