

Self-Supervised Graph Co-Training for Session-based Recommendation

Xin Xia

The University of Queensland
x.xia@uq.edu.au

Hongzhi Yin*

The University of Queensland
h.yin1@uq.edu.au

Junliang Yu

The University of Queensland
jl.yu@uq.edu.au

Yingxia Shao

BUPT
shaoyx@bupt.edu.cn

Lizhen Cui

Shandong University
clz@sdu.edu.cn

Problem Formulation

- **Session-based Recommendation:** Let $H = \{i_1, i_2, \dots, i_N\}$ denote the set of items, where N is the number of items. Each session is represented as a sequence $s = [i_{s,1}, i_{s,2}, \dots, i_{s,m}]$ ordered by timestamps and $i_{s,k}$ represents an interacted item of an anonymous user within the session s . Given s , the session-based recommendation task is to predict a ranked list $y = [y_1, y_2, \dots, y_N]$ where y_i ($1 \leq i \leq N$) is the corresponding predicted probability of item i .

Motivation

- The typical idea of applying self-supervised learning (SSL) to recommendation is conducting stochastic data augmentations by randomly dropping or masking some items/segments from the raw user-item interaction graph/sequence to create supervisory signals. However, this strategy may not be practicable for session-based recommendation, since the user interaction data generated in a session is much less than sequential recommendation.
- ***Self-supervised graph co-training:*** Two distinct graph (item view and session view) was constructed, first. Then, two asymmetric graph encoders were constructed over the two views, which recursively leverage the different connectivity information (internal and external) to generate ground-truth samples to supervise each other by contrastive learning.

Co-Training framework for session-based Recommendation

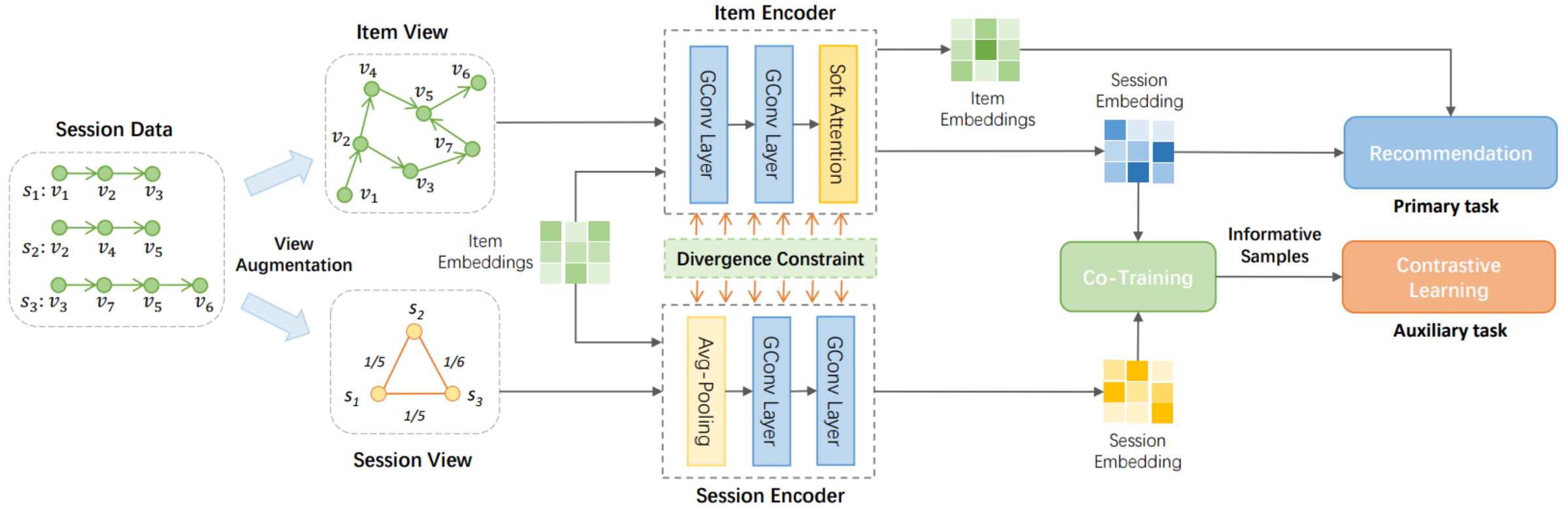


Figure 1: An overview of the proposed COTREC framework.

Item View Encoding

- **Information aggregation**

The item encoder with a simplified graph convolution layer for the item view is defined as

$$\mathbf{X}_I^{(l+1)} = \widehat{\mathbf{D}}_I^{-1} \widehat{\mathbf{A}}_I \mathbf{X}_I^{(l)} \mathbf{W}_I^{(l)}$$
$$\mathbf{x}_I^{t*} = \tanh(\mathbf{W}_1 [\mathbf{x}_I^t || \mathbf{p}_{m-t+1}] + \mathbf{b})$$

where \mathbf{x}_I^{t*} is the embedding of t -th item in session s . \mathbf{p}_{m-t+1} is a learnable position embedding. $\widehat{\mathbf{A}}_I = \mathbf{A}_I + \mathbf{I}$, $\widehat{\mathbf{D}}_{I,p,p} = \sum_{q=1}^m \widehat{\mathbf{D}}_{I,p,q}$, are the degree matrix and the adjacency matrix, \mathbf{I} is the identity matrix. $\mathbf{X}_I^{(l)}$ represent the l -th layer's item embeddings.

- **Session representation**

$$\theta_I = \sum_{t=1}^m \alpha_t \mathbf{x}_I^{t*}$$
$$\alpha_t = \mathbf{f}^T \sigma(\mathbf{W}_2 \mathbf{x}_s + \mathbf{W}_3 \mathbf{x}_I^{t*} + \mathbf{c}), \quad \mathbf{x}_s = \frac{1}{m} \sum_{t=1}^m \mathbf{x}_I^m$$

where \mathbf{x}_I^m is the item embedding within the session s after information aggregation. θ_I is the session embedding under item-view.

Session View Encoding

- **Information aggregation**

The session encoder is also a simplified graph convolution layer for the item view is defined as followed:

$$\Theta_S^{(l+1)} = \widehat{\mathbf{D}}_S^{-1} \widehat{\mathbf{A}}_S \Theta_S^{(l)} \mathbf{W}_S^{(l)}$$

$$\Theta_S^0 = \frac{1}{m} \sum_{t=1}^m \mathbf{x}_I^t || \mathbf{p}_{m-t+1}$$

where $\widehat{\mathbf{A}}_S = \mathbf{A}_S + \mathbf{I}$, $\widehat{\mathbf{D}}_S$ are the normalized adjacency matrix and degree matrix, respectively. $\Theta_S^{(l+1)}$ represents the l -th layer's session embeddings.

- **Session representation**

$$\Theta_S = \frac{1}{L+1} \sum_{l=0}^L \Theta_S^{(l)}$$

where Θ_S is the session embedding under session-view.

Co-Training

- **Positive/negative samples generation**

The positive/negative next-item samples generated via representation learned over the item and session views:

$$\mathbf{y}_I^p = \text{Softmax}(\text{score}_I^p), \text{score}_I^p = \mathbf{X}_I \boldsymbol{\theta}_I^p$$
$$\mathbf{y}_S^p = \text{Softmax}(\text{score}_S^p), \text{score}_S^p = \mathbf{X}^{(0)} \boldsymbol{\theta}_S^p$$

where $\mathbf{y}_I^p / \mathbf{y}_S^p$, $\boldsymbol{\theta}_I^p / \boldsymbol{\theta}_S^p$ are the predicted probability of and session embedding of session p , respectively, under item/session view. \mathbf{X}_I and $\mathbf{X}^{(0)}$ are item embedding.

$$\mathbf{c}_I^{p+} = \text{top} - K(\mathbf{y}_I^p), \mathbf{c}_S^{p+} = \text{top} - K(\mathbf{y}_S^p)$$

where $\mathbf{c}_I^{p+} / \mathbf{c}_S^{p+}$ are the positive samples under item/session view. $\text{top-K}(\cdot)$ means select items with the top-K highest confidence. For negative samples, we randomly select the samples from the items ranked in top 10% in \mathbf{y}_I^p excluding the positives to construct \mathbf{c}_I^{p-} .

Co-Training

- **Contrastive learning**

The InfoNCE was applied to maximize the lower bound of mutual information between the item pairs:

$$\begin{aligned} \mathcal{L}_{ssl} &= -\log \frac{\sum_{i \in c_I^{p+}} \psi(\mathbf{x}_I^{last}, \boldsymbol{\theta}_I^p, \mathbf{x}_I^i)}{\sum_{i \in c_I^{p+}} \psi(\mathbf{x}_I^{last}, \boldsymbol{\theta}_I^p, \mathbf{x}_I^i) + \sum_{j \in c_I^{p-}} \psi(\mathbf{x}_I^{last}, \boldsymbol{\theta}_I^p, \mathbf{x}_I^j)} \\ &\quad - \log \frac{\sum_{i \in c_S^{p+}} \psi(\mathbf{x}_{(0)}^{last}, \boldsymbol{\theta}_S^p, \mathbf{x}_{(0)}^i)}{\sum_{i \in c_S^{p+}} \psi(\mathbf{x}_{(0)}^{last}, \boldsymbol{\theta}_S^p, \mathbf{x}_{(0)}^i) + \sum_{j \in c_S^{p-}} \psi(\mathbf{x}_{(0)}^{last}, \boldsymbol{\theta}_S^p, \mathbf{x}_{(0)}^j)} \end{aligned}$$

where \mathbf{x}^{last} is the embedding of the last-clicked item of the given session. $\psi(x_1, x_2, x_3) = \exp\left(\frac{f(x_1 + x_2, x_3 + x_2)}{\tau}\right)$. $f(\cdot)$ is the cosine operation. τ is the temperature to amplify the effect of discrimination.

Objective Function

- **Divergence Constraint in Co-Training**

Due to the two encoders are training by the same data source, it somehow might lead to the model collapse problem. Therefore, adversarial examples are integrated into the training.

$$\mathcal{L}_{diff} = KL\left(Prob_I(\mathbf{X}_I), Prob_S(\mathbf{X}_I + \Delta_{adv}^I)\right) + KL\left(Prob_S(\mathbf{X}_I), Prob_I(\mathbf{X}_I + \Delta_{adv}^S)\right)$$

$$Prob_I(\mathbf{X}_I) = \text{softmax}(\mathbf{X}_I \boldsymbol{\theta}_I^p), Prob_S(\mathbf{X}_I) = \text{softmax}(\mathbf{X}_I \boldsymbol{\theta}_S^p)$$

$$\Delta_{adv} = \epsilon \frac{\Gamma}{||\Gamma||} \text{ where } \Gamma = \frac{\partial l_{adv}(\hat{y}|\mathbf{x}+\Delta)}{\partial \Delta}$$

The FGSM method was applied which adds adversarial perturbations on model parameters through fast gradient computation. ϵ is the control parameter.

$$\begin{aligned}\mathcal{L} &= \mathcal{L}_r + \beta \mathcal{L}_{ssl} + \alpha \mathcal{L}_{diff} \\ \mathcal{L}_r &= - \sum_{i=1}^N y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i) \\ \hat{\mathbf{z}}_i &= \theta_I^{S^T} \mathbf{x}_i, \hat{\mathbf{y}} = \text{softmax}(\hat{\mathbf{z}})\end{aligned}$$

Experiments

Method	Tmall				RetailRocket				Diginetica			
	P@10	M@10	P@20	M@20	P@10	M@10	P@20	M@20	P@10	M@10	P@20	M@20
FPMC	13.10	7.12	16.06	7.32	25.99	13.38	32.37	13.82	15.43	6.20	26.53	6.95
GRU4REC	9.47	5.78	10.93	---	---	---	---	---	---	---	---	---
NARM	19.17	10.42	23.30	---	---	---	---	---	---	---	---	---
STAMP	22.63	13.12	26.47	---	---	---	---	---	---	---	---	---
SR-GNN	23.41	13.45	27.57	---	---	---	---	---	---	---	---	---
GCE-GNN	28.01	15.08	33.42	---	---	---	---	---	---	---	---	---
S^2 -DHCN	26.22	14.60	31.42	---	---	---	---	---	---	---	---	---
COTREC	30.62	17.65	36.35	---	---	---	---	---	---	---	---	---

Table 2: Performance on Tmall

Method	Tmall	
	P@20	M@20
COTREC-base	27.71	13.36
base-DHCN	32.94	16.22
base-MASK	32.54	16.37
COTREC	36.35	18.04

Table 3: Comparisons of Different Components

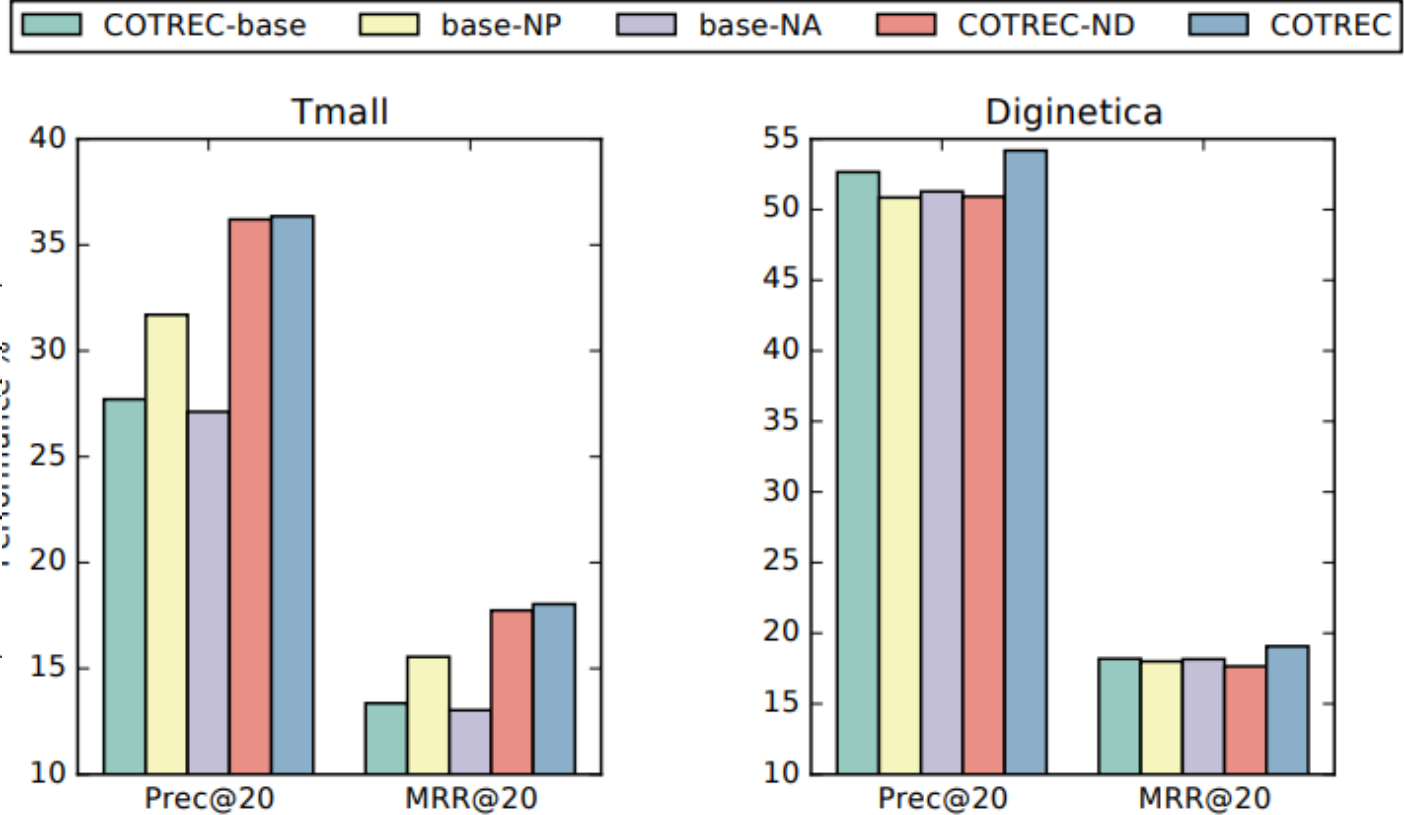


Figure 2: Ablation Study.

Experiments

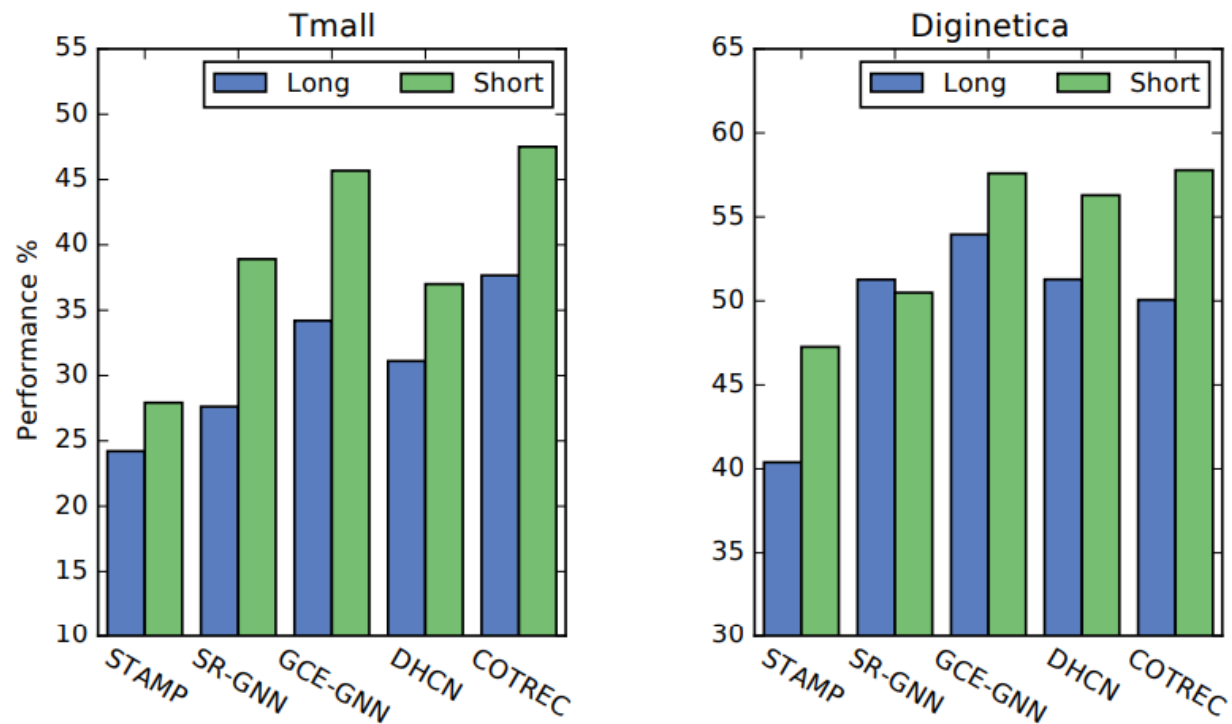


Figure 3: P@20 results on Long and Short.

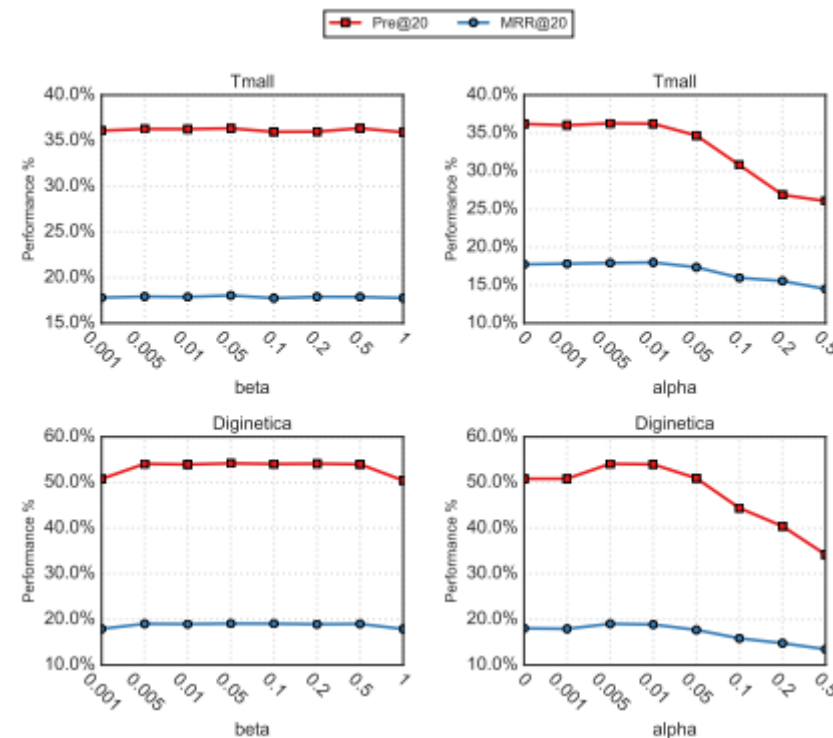


Figure 4: Hyperparameter Analysis.

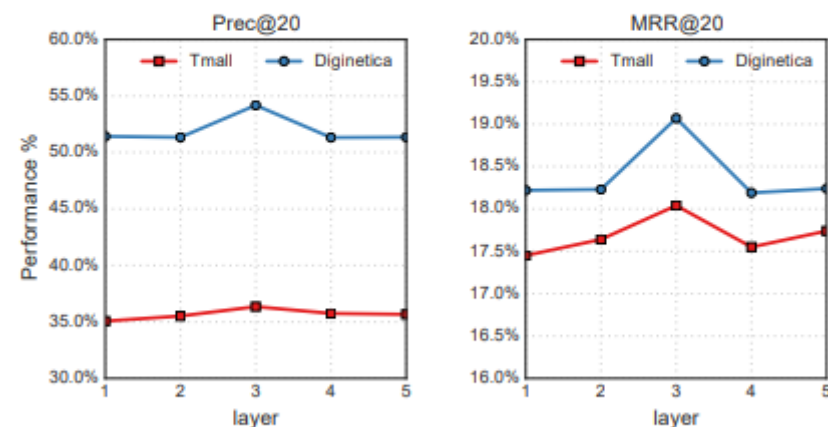


Figure 5: The impacts of the number of layer.