
DMSN final project: Improve LESSR model structure

TENG, LI-CHANG

Department of Electrical Engineering
National Cheng Kung University
n26091194@gs.ncku.edu.tw

TENG, LI-CHANG

Department of Electrical Engineering
National Cheng Kung University
n26091194@gs.ncku.edu.tw

TENG, LI-CHANG

Department of Electrical Engineering
National Cheng Kung University
n26091194@gs.ncku.edu.tw

TENG, LI-CHANG

Department of Electrical Engineering
National Cheng Kung University
n26091194@gs.ncku.edu.tw

Abstract

None

1 INTRODUCTION

Due to the highly practical value, session-based recommendation attracted researchers' great attention. Most of the methods proposed earlier are based on Markov chains or recurrent neural networks (RNNs). Recently, GNNs have become increasingly popular and achieved state-of-the-art performance in many tasks. There are also some attempts to apply GNNs to session-based recommendation.

Although these GNN-based methods obtained exciting results and offered a new and promising direction for session-based recommendation, we observe that there are two information loss problems in these methods. The first information loss problem in the existing GNN-based methods is called the lossy session encoding problem. The second information loss problem is called the ineffective long-range dependency capturing problem where these GNN-based methods cannot effectively capture all long-range dependencies.

To solve the above problems, author propose a novel GNN model called LESSR (Lossless Edge-order preserving aggregation and Shortcut graph attention for Session-based Recommendation). Through the structure of the original paper, we found that there are three places with GRU. We hope to increase the speed of the model by changing the part of GRU. We change GRU to attention, then change the model to Transformer, increase the number of layers and increase pos encoding, and finally see if the time and accuracy increase.

2 RELATED WORK

A given input session is first converted to a losslessly encoded graph called edge-order preserving (EOP) multigraph and a shortcut graph where the EOP multigraph could address the lossy session encoding problem and the shortcut graph could address the ineffective long-range dependency capturing problem. Then, the graphs along with the item embeddings are passed to multiple edge-order preserving aggregation (EOPA) and shortcut graph attention (SGAT) layers to generate latent

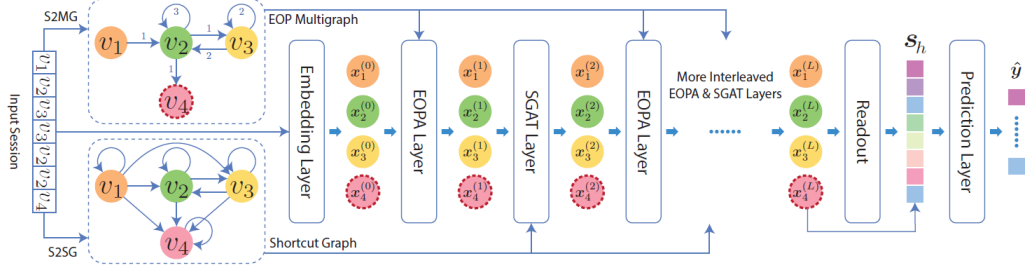


Figure 1: LESSR structure

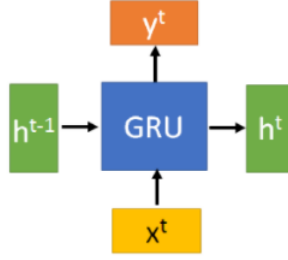


Figure 2: GRU structure

features of all nodes. The EOPA layers capture local context information using the EOP multigraph and the SGAT layers effectively capture long-range dependencies using the shortcut graph, and the architecture diagram is as figure 1.

3 PRELIMINARY

3.1 GRU

GRU (Gate Recurrent Unit) is a type of Recurrent Neural Network (RNN). It is proposed to solve the problems of long-term memory and gradients in back propagation. GRU input and output structure: The input-output structure of GRU is the same as that of ordinary RNN. There is a current input x^t and the hidden state h^{t-1} passed down from the previous node. This hidden state contains information about the previous node. The GRU structure as shown in figure 2.

3.2 Self-Attention

For self-attention, the three matrices Q(Query), K(Key), and V(Value) are all from the same input. First, we have to calculate the dot product between Q and K, and then in order to prevent the result from being too large, we will Divide by $\sqrt{d_k}$, then pass the basis activity function, and finally multiply by V, it can be expressed by the following formula 1.

$$Attention(Q, K, V) = softmax(\frac{QK^T}{\sqrt{d_k}})v \quad (1)$$

And multihead means that we can have different Q, K, V representations, and finally combine the results.

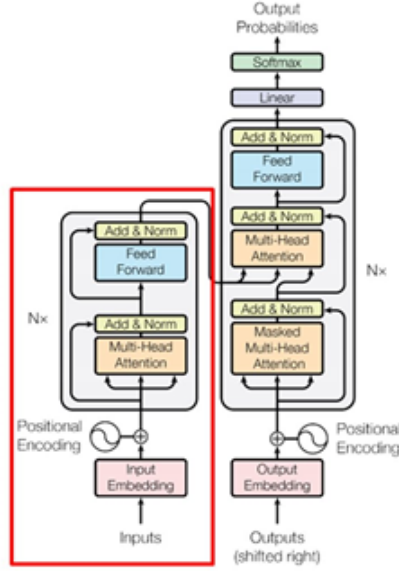


Figure 3: Transformer structure

4 METHODOLOGY

4.1 Transformer Encoder

In this section, we will introduce the structure of transformer encoder, which will be used to replace the GRU layer in EOPA in the following experiments.

Each transformer encoder has two sub-layers. The first one is Multi-Head Attention layer, and the second one is a fully connected feed forward layer. After both layer, we employ residual connection layers and layer normalization layers. We will specific each layer later. The Transformer structure is shown in figure 3.

4.2 Multi-Head Attention

The first one is the Multi-Head Attention layer. In this layer our goal is to get attention score from the graphs. In Multi-Head Attention layer, we have three groups of variables, Q, k, V, represent query, keys, value respectively. The number i in these groups (ex: $Q\{q_1, q_2, \dots, q_i\}$) is decided by how many heads we use. After that, we will use the following formula to compute the attention weight, where d_k represents vector dimensionality of k, v and Z represents final attention weight. The Multi-Head Attention structure is shown in figure 4.

4.3 Add & Norm

In these layers we use residual connection and layers normalization to help our model training. With residual connection, we can let gradient flow directive through the network, preventing gradient vanishing. As for layers normalization, it's effective at stabling the hidden state dynamic our model. In mathematical definition, residual connection can be simply denoted by the following formula 2.

$$add(x) = x + sublayer(x) \quad (2)$$

On the other hand, layers normalization is much complex. Assume we give input x over a mini batch size of m , which is $B = x_1, x_2, \dots, x_m$. Each sample x contain K elements. We first compute the mean and variance of each sample x , denoted by the following formula 3.

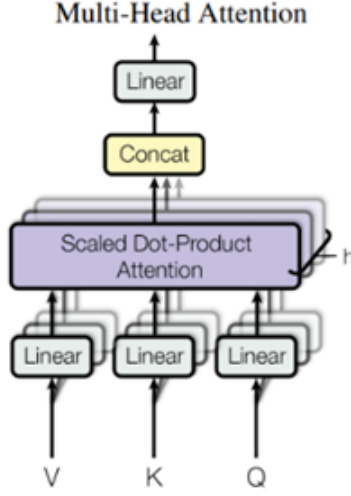


Figure 4: Multi-Head Attention

$$\mu_i = \frac{1}{K} \sum x_{i,k} \sigma_i^2 = \frac{1}{K} \sum (x_{i,k} - \mu_i)^2 \quad (3)$$

Then we normalize each sample such that the elements in the sample have zero mean and unit variance as formula 4. ϵ is for numerical stability in case the denominator becomes zero by chance.

$$x_{i,k} = \frac{(x_{i,k} - \mu_i)}{\sqrt{\sigma_i^2 - \epsilon}} \quad (4)$$

Combining these two layers we simplify its as formula 5

$$\text{LayerNorm}(x + \text{sublayer}(x)) \quad (5)$$

Recall that sublayers is Multi-Head Attention and feed forward network.

4.4 Feedforward Network

In addition to attend sub-layers, we employ a fully connected feed-forward network shown as formula 6 in our encoder. This consist of two linear transformation with a ReLU activation function.

$$\text{FFN}(x) = \max(0, xW_1 + b_1)W_2 + b_2 \quad (6)$$

5 EXPERIMENTS

In this section, we will introduce experiment setting, dataset, and analyze the experiment result. We conducted several experiments to check our hypotheses and evaluate our model with choosen metric.

5.1 Dataset

We choose Diginetica dataset¹ following LESSR [1] paper, which is the CIKM cup 2016 dataset provided by DIGINETICA Crop. There are 6 files in Diginetica dataset, but we only need the transaction one. As [1], we used last week sessions as test data. We got the same training and test

¹<https://competitions.codalab.org/competitions/111610>

Table 1: statistics of dataset

Diginetica	
No. of Clicks	981,620
No. of Sessions	777,029
No. of Items	42,596
Average length	4.80

Table 2: Multi-head w/o pos encoding

AGG.TYPE	HR@20	MRR@20	NDCG@20	Total impv.
baseline	52.82	18.3	25.93	-
Head=1	52.65	18.25	25.85	-0.903594775
Head=2	52.58	18.27	25.84	-0.965396084
Head=4	52.6	18.28	25.85	-0.834321462
Head=8	52.63	18.28	25.87	-0.700394057
Head=16	52.62	18.28	25.86	-0.757891648
Head=32	52.64	18.29	25.87	-0.626817026

set by following preprocessing method described in [1]. Statistics of Diginetica dataset is shown in Table 1.

5.2 Baseline and metrics

We choose [1] as our baseline model, then we tried to improve model structure in [1] by some changes. Comparing the metrics to [1], we could know the change is positive or negative influence. Following [1], the metrics we used are HR@20 (Hit Rate) and MRR@20 (Mean Reciprocal Rank).

5.3 Multi-Head Attention

MULTIHEADATTENTION² is an official implemented self-attention layer by pytorch. Here we replace GRU³ layer in EOPA block in [1] by MULTIHEADATTENTION layer. All settings are the same but GRU is now replaced by MULTIHEADATTENTION. We adjusted num of heads parameter in MULTIHEADATTENTION layer to see the influence of multi-head attention.

The pytorch official did not implement positional encoding in MULTIHEADATTENTION layer, so there is no position information within layer. To handle this problem we need to do position encoding manually. We found an official tutorial⁴ that manually implemented position encoding, so we followed the encoding method here.

Table 2 and Table 3 are the experiment result without and with positional encoding respectively. It turns out no matter multi-head or positional encoding, can not improve the result. So in next section we decided to use a more complex layer.

5.4 Transformer Encoder

In this section, we use transformer encoder⁵ to replace GRU. TransformerEncoder has a lot of hyperparameters, so we conducted 3 main experiments to tune the model: 1) dim_feedforward 2) nhead 3) encoder_layer. Also, each main experiment has two sub-experiments: 1) w/o pos encoding 2) with pos encoding.

²<https://pytorch.org/docs/stable/generated/torch.nn.MultiheadAttention.html>

³<https://pytorch.org/docs/stable/generated/torch.nn.GRU.html>

⁴https://pytorch.org/tutorials/beginner/transformer_tutorial.html

⁵<https://pytorch.org/docs/stable/generated/torch.nn.TransformerEncoder.html>

Table 3: Multi-head with pos encoding

AGG.TYPE	HR@20	MRR@20	NDCG@20	Total impv.
baseline	52.82	18.3	25.93	-
Head=1	52.57	18.26	25.83	-1.077538484
Head=2	52.55	18.29	25.85	-0.874337766
Head=4	52.59	18.3	25.88	-0.628267962
Head=8	52.62	18.29	25.87	-0.664681471
Head=16	52.54	18.31	25.86	-0.745415003
Head=32	52.57	18.32	25.89	-0.518277422

Table 4: dim exp w/o pos encoding

AGG.TYPE	HR@20	MRR@20	NDCG@20	Total impv.
baseline	52.82	18.3	25.93	-
Dim = 2048	52.73	18.25	25.83	-0.82926773
Dim = 1024	52.9	18.31	25.86	-0.063854988
Dim = 512	52.97	18.35	26	0.827164961
Dim = 256	52.67	18.28	25.88	-0.586099799
Dim = 128	52.88	18.34	25.94	0.370737939
Dim = 64	52.84	18.39	26.01	0.83819067
Dim = 32	52.7	18.24	25.85	-0.863578471
Dim = 16	52.68	18.3	25.88	-0.457877958

5.4.1 Dim_Feedforward Experiment

In this experiment we fix all hyperparameters but `dim_feedforward`. Table 4 shown the result without positional encoding. Table 5 shown the result with positional encoding.

From Table 4 and Table 5, we found that the best `dim_feedforward` is setting 512, whether with pos encoding or not, dimension 512 in both case has a good result, so we choose dimension 512 for our model in the later experiments.

5.4.2 Multi-Head Experiment

Here we fixed all hyperparameters but `nhead` to see the influence. Also, The `dim_feedforward` set to 512. Result without positional encoding shown in Table 6 and result with positional encoding shown in Table 7.

Comparing Table 6 and Table 7, We found the metrics without positional encoding are usually better than the other one. So positional information might not a critical info in this scenario.

Note that best performance appeared when `nhead` set to 8 with positional encoding.

Table 5: dim exp with pos encoding

AGG.TYPE	HR@20	MRR@20	NDCG@20	Total impv.
baseline	52.82	18.3	25.93	-
Dim = 2048	52.74	18.28	25.88	-0.45357424
Dim = 1024	52.86	18.3	25.88	-0.117097951
Dim = 512	52.85	18.36	25.97	0.538926994
Dim = 256	52.7	18.26	25.86	-0.715723485
Dim = 128	52.89	18.37	25.95	0.592169956
Dim = 64	52.79	18.34	25.95	0.238913304
Dim = 32	52.74	18.29	25.9	-0.321798695
Dim = 16	52.6	18.36	25.92	-0.127205414

Table 6: multi-head exp w/o pos encoding

AGG.TYPE	HR@20	MRR@20	NDCG@20	Total impv.
baseline	52.82	18.3	25.93	-
nhead=1	52.97	18.35	26	0.827164961
nhead=2	52.77	18.37	25.95	0.364983285
nhead=4	52.98	18.35	26	0.846097184
nhead=8	52.87	18.37	25.96	0.592870879
nhead=16	52.78	18.37	25.97	0.461046244
nhead=32	52.92	18.41	25.97	0.944676596

Table 7: multi-head exp with pos encoding

AGG.TYPE	HR@20	MRR@20	NDCG@20	Total impv.
baseline	52.82	18.3	25.93	-
nhead=1	52.85	18.36	25.97	0.538926994
nhead=2	52.7	18.31	25.91	-0.2496726
nhead=4	52.88	18.38	25.98	0.743578647
nhead=8	52.87	18.41	26.03	1.081407692
nhead=16	52.82	18.35	25.96	0.388920149
nhead=32	52.75	18.35	25.95	0.217829222

5.4.3 Num-Layers Experiment

Here all hyperparameters was fixed but num_layers will be change. The dim_feedforward was set to 512 and nhead was set to 1. Table 8 and Table 9 shown the result without and with positional encoding respectively.

We found that whether transformer encoder with positional encoding or not, it has similar trend that performance decreased as layers increased. And we found as layers increased, models's inference time also increased.

6 CONCLUSION

In our experimets, no matter with postional encoding or not, the performance of MUTIHEADATTENTION is worse than GRU, although using MUTIHEADATTENTION is slightly fast. In muti-head experiment of TransformerEncoder, we found there is no distinct trend about num of head, also we found when TransformerEncoder stacking more layers, model's performance will decreased, this means the model is over-fitting, moreover, evaluate time will increased, this cause training time getting longer.

After several experimets, we found that TransformerEncoder can improve model performance by replacing GRU in EOPA layer, but the parameter of TransformerEncoder need in proper setting. In SGAT and Readout layer, we tried to change the attention mechanisms, but we didn't got a good

Table 8: num-layers exp w/o pos

AGG.TYPE	HR@20	MRR@20	NDCG@20	Total impv.
baseline	52.82	18.3	25.93	-
layer=1	52.97	18.35	26	0.827164961
layer=2	52.72	18.41	25.93	0.41177067
layer=3	52.69	18.4	25.95	0.37745993
layer=4	52.69	18.33	25.89	-0.236445941
layer=6	52.65	18.13	25.72	-2.060682268
layer=8	52.24	17.96	25.48	-4.691433984
layer=16	52.11	17.77	25.34	-6.515719401

Table 9: num-layers exp with pos

AGG.TYPE	HR@20	MRR@20	NDCG@20	Total impv.
baseline	52.82	18.3	25.93	-
layer=1	52.85	18.36	25.97	0.538926994
layer=2	52.77	18.41	25.96	0.622127888
layer=3	52.72	18.34	25.88	-0.163569833
layer=4	52.84	18.32	25.9	0.031457958
layer=6	52.8	18.18	25.78	-1.272082675
layer=8	52.3	17.95	25.46	-4.709616194
layer=16	52.04	17.81	25.37	-6.313969619

result. This might because attention mechanisms is dependent on model structure and data, it won't be useful if we just replace attention mechanisms from different paper.

Finally the best result we got is when layer=1 and nhead=8 with positional encoding, total improvement of three metrics is 1.08%, although there is still a space for improvements in this model. The result is acceptable to us because we didn't change the model structure so much but just change a layer inside EOPA block. If we could modify SGAT and Readout layer's attention mechanisms to multi-head attention, we might get a better result.

References

- [1] Tianwen Chen and Raymond Chi-Wing Wong. Handling information loss of graph neural networks for session-based recommendation. In *Proceedings of the 26th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD '20)*, pages 1172—1180, 2020.