# Neural Common Neighbor with Completion for Link Prediction

## ICLR 2024

Xiyuan Wang, Haotong Yang, Muhan Zhang

Institute of Artificial Intelligence, Peking University

April 25, 2024
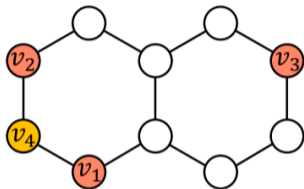
# MPNN for Link Prediction

Xiyuan Wang,
Haotong Yang,
Muhan Zhang

Neural Common
Neighbor

Completion for
Input Graph

Experiments

Vanilla MPNN fails in this task



- Learns node representation only.
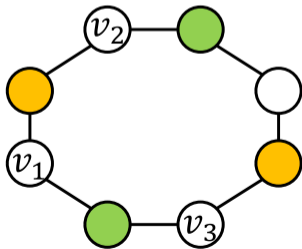- Cannot distinguish $(v_1, v_2)$ and $(v_1, v_3)$.

Structural feature like number of common neighbor can help.

# Structure Feature cannot Capture Node Feature

Xiyuan Wang,
Haotong Yang,
Muhan Zhang

Commonly structure features

$$
\begin{array}{ll}
\text{Common Neighbor} & \sum_{u\in N(i)\cap N(j)} 1 \\
\text{Resource Allocation} & \sum_{u\in N(i)\cap N(j)} \frac{1}{d(u)} \\
\text{Adamic Adar} & \sum_{u\in N(i)\cap N(j)} \frac{1}{\log d(u)}
\end{array}
\tag{1}
$$

- unlearnable
- unable to capture node feature



Cannot distinguish $(v_1, v_2)$ and $(v_1, v_3)$.

# New Architecture

Xiyuan Wang,
Haotong Yang,
Muhan Zhang

Neural Common Neighbor

Completion for Input Graph

Experiments

Structural features (SF) like common neighbor are widely used.

- Existing works utilize SF in two ways.
  - SF-then-MPNN. Take SF as MPNN's input
    - Low scalability, need to rerun MPNN as the SF changes with target link.
  - SF-and-MPNN. Ensemble MPNN with SF.
    - MPNN and SF are completely separated. Low expressivity.
- We use SF to guide the pooling of MPNN's output (MPNN-then-SF).
  - Good scalability and expressivity.

Xiyuan Wang,
Haotong Yang,
Muhan Zhang

$$\left. \begin{array}{ll} \text{Common Neighbor} & \sum_{u \in N(i) \cap N(j)} 1 \\ \text{Resource Allocation} & \sum_{u \in N(i) \cap N(j)} \frac{1}{d(u)} \\ \text{Adamic Adar} & \sum_{u \in N(i) \cap N(j)} \frac{1}{\log d(u)} \\ \text{Neo-GNN} & \sum_{u \in N_1^{l_1}(i) \cap N_1^{l_2}(j)} A_{iu}^{l_1} A_{ju}^{l_2} f(d(u)) \\ \text{BUDDY} & \sum_{u \in N_{l_1}^1(i) \cap N_{l_2}^1(j)} 1, \sum_{u \in N_l^1(i) - \bigcup_{l'=1}^k N_{l'}^1(j)} 1 \end{array} \right\} \Rightarrow \quad (2)$$

$$\Rightarrow \sum_{u \in N_{l_1}^{l_2}(i) \oplus N_{l_1'}^{l_2'}(j)} g(A_{iu}^{l_2}) g(A_{ju}^{l_2'}) \text{MPNN}(u, A, X) \quad (3)$$

Higher-order neighbors and neighborhood differences lead to negligible performance gain, leading to our NCN model:

$$\text{NCN}(i, j, A, X) = \sum_{u \in N(i) \cap N(j)} \text{MPNN}(u, A, X) \quad (4)$$

Xiyuan Wang,
Haotong Yang,
Muhan Zhang

Incompleteness of graph is ubiquitous in link prediction tasks.

- The target edge exists in input graph on training set but not on test set.

Besides the target edge, other edges, like those connected to common neighbors, is also affected.

To visualize it, we assume that

- Graph with training set edges only is the *incomplete* graph.
- Graph with training, validation and test set edges is the *complete* graph.



Figure: Ogbl-collab dataset (a) distribution of common neighbor (b) performance of common neighbor

PEKING UNIVERSITY

Xiyuan Wang,
Haotong Yang,
Muhan Zhang

Neural Common
Neighbor

Completion for
Input Graph

Experiments



(a)    (b)

- blue and green lines in (a): a significant *distribution shift* between the training and test sets in the incomplete graph of the ogbl-collab dataset.
- red and orange lines: shift disappears when the graph is complete.

Distribution shifts can enlarge the gap between training and test error.

Xiyuan Wang,
Haotong Yang,
Muhan Zhang

Neural Common
Neighbor

Completion for
Input Graph

Experiments



(a)  (b)

- Blue and green lines in (a): there are fewer common neighbors in the incomplete graph.

Loss of Common Neighbor Information can lead to high training error and thus high test error.

Xiyuan Wang,
Haotong Yang,
Muhan Zhang

Neural Common Neighbor

Completion for Input Graph

Experiments



(a)

(b)

- (b): Imcompleteness of non-target links leads to significant performance degradation.

Xiyuan Wang,
Haotong Yang,
Muhan Zhang

Due to imcompleteness, NCN can use not only the common neighbors in the input graph. It can also predict unobserved common neighbors.

Given a target link $(i, j)$, the probability that $u$ is a common neighbor of $(i, j)$ is

$$
P_{uij} = \begin{cases} 1 & \text{if } u \in N(i,A) \cap N(j,A) \\ \sigma(\text{NCN}(i, u, A, X)) & \text{if } u \in N(j,A) - N(i,A) \\ \sigma(\text{NCN}(j, u, A, X)) & \text{if } u \in N(i,A) - N(j,A) \\ 0 & \text{otherwise} \end{cases} \tag{5}
$$

NCN with Completion (NCNC) becomes

$$
\text{NCNC}(i, j, A, X) = \sum_{u \in V} P_{uij} \text{MPNN}(u, A, X). \tag{6}
$$

# Link Prediction

Xiyuan Wang,
Haotong Yang,
Muhan Zhang

Neural Common Neighbor

Completion for Input Graph

Experiments

| | Cora | Citeseer | Pubmed | Collab | PPA | Citation2 | DDI |
|---|---|---|---|---|---|---|---|
| Metric | HR@100 | HR@100 | HR@100 | HR@50 | HR@100 | MRR | HR@20 |
| **CN** | $33.92_{\pm0.46}$ | $29.79_{\pm0.90}$ | $23.13_{\pm0.15}$ | $56.44_{\pm0.00}$ | $27.65_{\pm0.00}$ | $51.47_{\pm0.00}$ | $17.73_{\pm0.00}$ |
| **AA** | $39.85_{\pm1.34}$ | $35.19_{\pm1.33}$ | $27.38_{\pm0.11}$ | $64.35_{\pm0.00}$ | $32.45_{\pm0.00}$ | $51.89_{\pm0.00}$ | $18.61_{\pm0.00}$ |
| **RA** | $41.07_{\pm0.48}$ | $33.56_{\pm0.17}$ | $27.03_{\pm0.35}$ | $64.00_{\pm0.00}$ | $49.33_{\pm0.00}$ | $51.98_{\pm0.00}$ | $27.60_{\pm0.00}$ |
| **GCN** | $66.79_{\pm1.65}$ | $67.08_{\pm2.94}$ | $53.02_{\pm1.39}$ | $44.75_{\pm1.07}$ | $18.67_{\pm1.32}$ | $84.74_{\pm0.21}$ | $37.07_{\pm5.07}$ |
| **SAGE** | $55.02_{\pm4.03}$ | $57.01_{\pm3.74}$ | $39.66_{\pm0.72}$ | $48.10_{\pm0.81}$ | $16.55_{\pm2.40}$ | $82.60_{\pm0.36}$ | $53.90_{\pm4.74}$ |
| **SEAL** | $81.71_{\pm1.30}$ | $83.89_{\pm2.15}$ | $75.54_{\pm1.32}$ | $64.74_{\pm0.43}$ | $48.80_{\pm3.16}$ | $87.67_{\pm0.32}$ | $30.56_{\pm3.86}$ |
| **NBFnet** | $71.65_{\pm2.27}$ | $74.07_{\pm1.75}$ | $58.73_{\pm1.99}$ | OOM | OOM | OOM | $4.00_{\pm0.58}$ |
| **NeoGNN** | $80.42_{\pm1.31}$ | $84.67_{\pm2.16}$ | $73.93_{\pm1.19}$ | $57.52_{\pm0.37}$ | $49.13_{\pm0.60}$ | $87.26_{\pm0.84}$ | $63.57_{\pm3.52}$ |
| **BUDDY** | $88.00_{\pm0.44}$ | $\underline{92.93_{\pm0.27}}$ | $74.10_{\pm0.78}$ | $\underline{65.94_{\pm0.58}}$ | $49.85_{\pm0.20}$ | $87.56_{\pm0.11}$ | $78.51_{\pm1.36}$ |
| **NCN** | $\underline{89.05_{\pm0.96}}$ | $91.56_{\pm1.43}$ | $\underline{79.05_{\pm1.16}}$ | $64.76_{\pm0.87}$ | $\underline{61.19_{\pm0.85}}$ | $\underline{88.09_{\pm0.06}}$ | $\underline{82.32_{\pm6.10}}$ |
| **NCNC** | $\mathbf{89.65_{\pm1.36}}$ | $\mathbf{93.47_{\pm0.95}}$ | $\mathbf{81.29_{\pm0.95}}$ | $\mathbf{66.61_{\pm0.71}}$ | $\mathbf{61.42_{\pm0.73}}$ | $\mathbf{89.12_{\pm0.40}}$ | $\mathbf{84.11_{\pm3.67}}$ |

# Ablation Study

Xiyuan Wang,
Haotong Yang,
Muhan Zhang

Neural Common Neighbor

Completion for Input Graph

Experiments

| | Cora | Citeseer | Pubmed | Collab | PPA | Citation2 | DDI |
|---|---|---|---|---|---|---|---|
| Metric | HR@100 | HR@100 | HR@100 | HR@50 | HR@100 | MRR | HR@20 |
| **CN** | $33.92_{\pm0.46}$ | $29.79_{\pm0.90}$ | $23.13_{\pm0.15}$ | $56.44_{\pm0.00}$ | $27.65_{\pm0.00}$ | $51.47_{\pm0.00}$ | $17.73_{\pm0.00}$ |
| **GAE** | $89.01_{\pm1.32}$ | $91.78_{\pm0.94}$ | $78.81_{\pm1.64}$ | $36.96_{\pm0.95}$ | $19.49_{\pm0.75}$ | $79.95_{\pm0.09}$ | $61.53_{\pm9.59}$ |
| **GAE+CN** | $88.61_{\pm1.31}$ | $91.75_{\pm0.98}$ | $79.04_{\pm0.83}$ | $64.47_{\pm0.14}$ | $51.83_{\pm0.58}$ | $87.81_{\pm0.06}$ | $80.71_{\pm5.56}$ |
| **NCN2** | $88.87_{\pm1.34}$ | $91.36_{\pm1.02}$ | $80.21_{\pm0.78}$ | $65.43_{\pm0.46}$ | OOM | OOM | OOM |
| **NCN-diff** | $89.12_{\pm1.04}$ | $91.96_{\pm1.23}$ | $80.28_{\pm0.88}$ | $64.08_{\pm0.40}$ | $57.86_{\pm1.26}$ | $86.68_{\pm0.16}$ | $17.67_{\pm8.70}$ |
| **NCN** | $89.05_{\pm0.96}$ | $91.56_{\pm1.43}$ | $79.05_{\pm1.16}$ | $64.76_{\pm0.87}$ | $61.19_{\pm0.85}$ | $88.09_{\pm0.06}$ | $82.32_{\pm6.10}$ |
| **NoTLR** | $85.46_{\pm1.65}$ | $88.08_{\pm1.23}$ | $76.59_{\pm1.33}$ | $64.22_{\pm0.49}$ | $60.66_{\pm0.63}$ | $88.64_{\pm0.14}$ | $66.52_{\pm11.37}$ |
| **NCNC** | $89.65_{\pm1.36}$ | $93.47_{\pm0.95}$ | $81.29_{\pm0.95}$ | $66.61_{\pm0.71}$ | $61.42_{\pm0.73}$ | $89.12_{\pm0.40}$ | $84.11_{\pm3.67}$ |
| **NCNC-2** | $89.14_{\pm0.84}$ | $93.14_{\pm0.96}$ | $81.41_{\pm1.07}$ | $66.80_{\pm0.43}$ | $>24h$ | $>24h$ | $>24h$ |

# Scalability Comparison

Xiyuan Wang,
Haotong Yang,
Muhan Zhang

Neural Common Neighbor

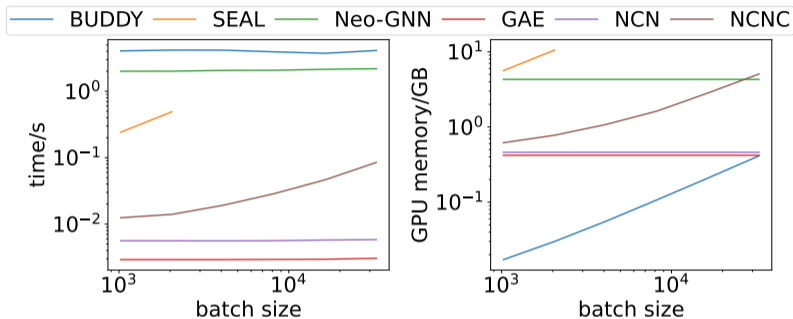Completion for Input Graph

Experiments



Figure: Inference time and GPU memory on ogbl-collab. The process we measure includes preprocessing and predicting one batch of test links. Relation between time $y$ and batch size $t$ is $y = B + Ct$, where $B, C$ are model specific constants. SEAL has out-of-memory problem and only uses small batch sizes.