

ニューラルネットワークを用いた異種データのグラフ埋め込み

奥野 彰文^{1,2}, 下平 英寿^{1,2}

¹ 京都大学大学院 情報学研究科 ² 理化学研究所 革新知能統合研究センター

画像やテキストといった, $D(\geq 2)$ 種類のデータを異種データと呼ぶ. 各種類 $d = 1, 2, \dots, D$ での i 番目のデータを表すベクトル $\mathbf{x}_i^{(d)} \in \mathbb{R}^{p_d}$ をデータベクトルと呼ぶ. $\mathbf{x}_i^{(d)}, \mathbf{x}_j^{(e)}$ 間の関連の強さを表す重み $w_{ij}^{(de)} = w_{ji}^{(ed)} \geq 0$ をマッチングウェイトと呼ぶ. $\{\mathbf{x}_i^{(d)}\}_{i=1}^{n_d}$ と $\{w_{ij}^{(de)}\}_{i,j=1}^{n_d, n_e}$ は考慮するすべての (d, e) についての所与の観測値であり, データベクトルの数 n_d と次元 p_d はデータの種類 d に応じて違ってよい.

次元が種類 d に応じて異なるデータベクトルを統計解析することは, 一般に難しい. そこで, $\{\mathbf{x}_i^{(d)}\}, \{w_{ij}^{(de)}\}$ を考慮して連続変換 $f_\psi^{(d)}: \mathbb{R}^{p_d} \rightarrow \mathbb{R}^K$ を学習し, $\mathbf{x}_i^{(d)} \in \mathbb{R}^{p_d}$ を次元が d に依存しない

$$\mathbf{y}_i^{(d)} := f_\psi^{(d)}(\mathbf{x}_i^{(d)}) \in \mathbb{R}^K \quad (1)$$

に変換することで, 既存の最近傍探索やクラスタリングの手法を異種データにも適用できる.

しかしながら, 既存の異種データ解析手法はすべて, (i) 一対一対応を仮定するなどマッチングウェイト $w_{ij}^{(de)}$ に制約が必要, (ii) $f_\psi^{(d)}$ に線形モデルを仮定し表現力が乏しい, (iii) 計算量が大きい, のうち少なくとも一つの欠点を克服できなかった. そこで, 本研究ではポアソン分布を用いた確率モデル

$$\begin{aligned} w_{ij}^{(de)} \mid \mathbf{x}_i^{(d)}, \mathbf{x}_j^{(e)} &\stackrel{\text{indep.}}{\sim} \text{Po} \left(\exp \left(h(\mathbf{x}_i^{(d)}, \mathbf{x}_j^{(e)}; \boldsymbol{\theta}) \right) \right), \\ h(\mathbf{x}_i^{(d)}, \mathbf{x}_j^{(e)}; \boldsymbol{\theta}) &:= \langle \mathbf{y}_i^{(d)}, \mathbf{y}_j^{(e)} \rangle - \gamma^{(de)}, \quad (\boldsymbol{\theta} = (\boldsymbol{\psi}, \{\gamma^{(de)}\}_{d,e})), \end{aligned} \quad (2)$$

($\langle \cdot, \cdot \rangle$ は内積) を基にして, 尤度を minibatch SGD により最大化することで変換 (1) のパラメータ $\boldsymbol{\psi}$ を求める枠組み Probabilistic Multi-view Graph Embedding (PMvGE) を提案する. PMvGE は既存手法の問題点 (i)–(iii) をすべて同時に解決する. $f_\psi^{(d)}$ には任意の連続変換が利用できるが, 本研究では特にニューラルネットワークと線形変換の2つを用いた. それぞれの変換について更に以下の結果を得た.

- (A) ニューラルネットワーク (NN). 実数値 NN $u_\xi^{(d)}: \mathbb{R}^{p_d} \rightarrow \mathbb{R}$ による変換 $r_i^{(d)} := u_\xi^{(d)}(\mathbf{x}_i^{(d)})$ を用いたモデル $\langle \mathbf{y}_i^{(d)}, \mathbf{y}_j^{(e)} \rangle + r_i^{(d)} + r_j^{(e)}$ を Shifted Inner Product Similarity (SIPS) と呼び, その特殊形である (2) を特に Constantly-SIPS (C-SIPS) と呼ぶ. SIPS と C-SIPS はベクトル値 NN の内積を利用しており, 実数値 NN の表現定理 (Cybenko, 1989) を適用できない. 本研究では, まず SIPS が単なる内積類似度よりも表現能力が高く, 任意の条件付き正定値類似度を任意の精度で近似できることを示した. この結果により, SIPS を用いたグラフ埋め込みが漸近的にはポアンカレ埋め込み (Nickel and Kiela, 2017) などと同等の表現能力を持つことが分かった. さらに, C-SIPS を用いた場合についても, $\gamma^{(de)} \geq 0$ を十分に大きくすれば, SIPS と同等の表現能力があることを示した.
- (B) 線形変換. 線形変換を利用するとき, PMvGE の近似解が Cross-Domain Matching Correlation Analysis (CDMCA) (Shimodaira, 2016) と等しくなることを示した. CDMCA は HIMFAC (Nori et al., 2012), CvGE (Huang et al., 2013), CCA, PCA などの一般化であるから, PMvGE はこれらの様々な多変量解析手法を非線形に拡張し, 計算を高速化したより一般的な枠組みといえる.

実データを用いた実験でも, PMvGE が Deep CCA などの既存手法の性能を上回った.

なお, 本講演は以下の2つの発表を基にしている.

Okuno, A., Hada, T., and Shimodaira, H. (2018). A probabilistic framework for multi-view feature learning with many-to-many associations via neural networks. In *Proceedings of the International Conference on Machine Learning (ICML)*.

Okuno, A. and Shimodaira, H. (2018). On representation power of neural network-based graph embedding and beyond. In *Proceedings of the Theoretical Foundations and Applications of Deep Generative Models workshop at ICML*.