# 半教師有りグラフニューラルネットワークを用いた CRUD関係に基づくシステム移行単位の最適化

Optimization for the Migration Units in Information Systems based on CRUD Relations with Semi-Supervised Graph Neural Networks

## 貞光 九月 \*1 Kugatsu Sadamitsu

\*1フューチャー株式会社 Strategic AI Group

SAIG, Future Corporation

Migrating large-scale information system as a whole is impossible due to its size, thus it is necessary to divide the system into manageable migration units. Typically, optimizing the migration units can only be done by experienced engineers, and takes large amount of time. In this paper, we propose a method finding optimal migration units. To find optimal migration units, we interpret CRUD relations between functions and tables as a graph structure, and adapt semi-supervised graph neural network, Diffusion Convolutional Neural Networks, with a graph embedding method. Furthermore, we unify CRUD relation with file arrangement information and attributes of each node. In experiments, the proposed methods find practical migration units that can be applied to real examples.

# 1. はじめに

産業の基幹となる大規模な情報システムにおいて,外部環 境変化への対応や業務の最適化等の目的から,システム刷新の 要求は必然的に生じる [木原他 2015].大規模な情報システム の刷新を行う際,一度にすべての機能を新システムに切り替え ることはリスクが高く非現実的であり,実際の開発では,一部 の機能やテーブルの集合を部分的に移行単位として切り出し, 段階的に移行するという安全策が採用されることが多い.この 時の移行単位の設計は,業務とプログラムに精通する技術者が 自身の経験則に基づいて決定することが多く,設計自体に多大 な工数を要するという課題や,対象システムに対する知識と経 験を持つ一部の技術者にスキルが偏在してしまうという属人性 の課題があった.我々はこれらの課題を解決するために,大規 模システムの段階移行開発における,移行単位の自動最適化手 法を提案する.

本稿では,移行単位最適化の手がかりとして,システム設計 時に用いられる,機能とテーブルの操作関係を記述した CRUD 表に着目する. CRUD とは、テーブルに対する代表的な操作 である, Create, Read, Update, Delete の頭文字をとったも のであり、機能からテーブルへの操作がどのように行われるか 記述したものである. CRUD 表は人手による移行単位設計の 際にも参照されることが多く、自動化においても有益なリソー スとして活用できることが見込まれる.また、2.節で述べる 通り、CRUD 表は一種のグラフ構造とみなすことができ、そ のためシステム移行単位の最適化問題は、CRUD 関係を表す グラフにおける各機能ノード及び各テーブルノードの分類問 題と定義できる. さらに本稿では CRUD 関係に加え,機能の ソースコードやテーブルデータのコンピュータ内のディレクト リ配置情報や、機能やテーブル自体が持つ属性情報を含めた統 一的な1つのグラフとして取り扱うことで、より高精度な解析 を可能とする.

グラフ中のノード分類手法はこれまでに多くの手法 が提案されており、本稿では半教師有りグラフノード分 類手法の一種で、近年高い精度が報告されているグラフ ニューラルネットワークの一種, Diffusion Convolutional Neural Network [Atwood and Towsley 2016] を用い, さ らにグラフ埋め込み (graph embedding) 法の一種である node2vec [Grover and Leskovec 2016] と Poincaré embedding [Nickel and Kiela 2017] をそれぞれ併用することで,比 較的少量の教師データのみを用いた場合でも,高い精度で移行 単位の推定が可能となることを示す.

# 移行単位最適化のための CRUD 関係の利 用

# 2.1 CRUD 関係と CRUD 表

本節では、移行単位の最適化に必要な手がかりとして、CRUD 関係および CRUD 表を用いることを考える.システム開発に おいて、機能とテーブルの関係は極めて重要である.システム の設計や開発の各フェーズにおいて、機能とテーブルの操作関 係を漏れなく設計・開発に落とし込むために広く用いられるも のが CRUD 表である.図1 (a) に CRUD 表の例を、図1 (b) に CRUD 表に対応するシステム中の CRUD 関係を表す.例え ば、図1(a) において、機能ノード "Func2"は、テーブルノー ド "Table2" に対し "Create"の関係を持ち、"Table3" に対 し "Read"の関係を持つことを表す. CRUD 関係を持つノー ド同士は比較的強い結びつきを持つと考えられることから、本 稿ではこの CRUD 表に着目した移行単位の最適化を行う.

本稿で扱うタスク概要についても同図を用いて表す.まず図 1 (a) 中, "Unit class1", "Unit class2"と記した囲いが, 人 手で解析した移行単位の理想的な分類結果であるとする.本稿 で扱う移行単位最適化は,一部の機能またはテーブルに対応す る移行単位クラスを教師信号として与え (図 (b) 中 "seed"と 記した "Func2"と "Func5"), その他のノードの移行単位クラ スを正しく分類することが目的となる.

さらに CRUD 表は, CRUD 関係をグラフとみなした場合 の隣接行列とみなすことが可能であり,実際に,図1 (b)の CRUD 関係に対応するグラフ構造を,図1(a)の CRUD 表と いう隣接行列に等価に変換することができる.システム操作関 係をグラフとして解釈することにより,3.節で述べるように, グラフ解析のための機械学習を直接的に適用できるという利点 が生じる.

連絡先: 貞光九月,フューチャー株式会社,東京都品川区大
崎1-2-2アートヴィレッジ大崎セントラルタワー, k.sadamitsu.ic@future.co.jp



図 1: (a)CRUD 表と人手付与した正解移行単位 (Unit class 1, 2), (b) CRUD 表に対応する機能とテーブル間の関係図. "seed" と記す "Func2" と "Func5" が教師データで,残る全 ノードがクラス推定対象ノード. (c) 機能とテーブル間の関係 図に,ファイル配置情報を重畳したもの. Func1 と Func2 が 同じクラスに属す可能性を示している.

# 2.2 ファイル配置情報とノード属性の利用

CRUD 関係は機能とテーブルの関係を表現する上で重要な 手がかりとなりえるが、本節では CRUD 関係以外で追加でき る手がかり情報を検討する.

1つ目の手がかりは、各機能のソースコードやテーブルデー タのコンピュータ上でのファイル配置情報である.システム開 発を行う上で、同じディレクトリに配置される機能やテーブル は、1つの単位としてまとめて移行する方が適切であるとい う仮説に基づいて導入する.ファイル配置情報の利用方法に はいくつか方策が考えられるが、本稿では直接 CRUD 関係の グラフ構造に重畳する方法を提案する.具体的には,元とな る CRUD 関係を表すグラフ構造に対し、OS のディレクトリ に相当する仮想ノードを設け、そのディレクトリ内に配置され る機能のソースコード・テーブルデータ,あるいは親・子ディ レクトリに対し、CRUD 関係と同様にエッジを張ることとす る. 実際に重畳した具体例を図1の(c)に示す. このような重 畳を行うことで,図1の(b)における "Func1"と "Table1" のように他ノードから孤立したサブグラフが、(c)ではファイ ル配置情報を考慮することで "Func2" 等と同じディレクトリ "/a/"以下の配置であることが分かり, "Func2"と同じ移行 単位とすべきと判断する根拠とすることができる. なおグラ フ中のエッジを定義する際に、CRUD 関係の種別や、ファイ ル配置情報か否かといったエッジの種類を複数分けて想定した り、エッジの重みとして考慮することも考えられるが、本稿で は全てのエッジを同一種、同一重みのエッジとして扱うことと した.

2つ目の手がかりは、ノード自体に付与される属性情報であ る. CRUD 関係やファイル配置情報は機能・テーブルノード 間のエッジに関する情報であるが、ノード自体にも情報が付属 する場合がある.例えば、テーブルの場合「マスタ」「トラン ザクション」といったテーブルの種類を表す属性が付与されて いる場合が多くある.これらの属性を考慮することで、重要な 属性値を持つ機能ノードとテーブルノード間に CRUD 関係が ある場合,他の CRUD 関係よりも優先してそれらノードを同 じクラスにまとめるべきである、といった指針を得ることが可 能となる.次節で述べるグラフノード分類法の中には、各ノー ドの属性を直接取り入れることが可能な手法が多く存在し、こ れら属性情報を導入することは容易である.

# 3. 半教師有りグラフニューラルネットワーク を用いた移行単位最適化

#### 3.1 半教師有りグラフニューラルネットワーク

前節において, CRUD 関係を一種のグラフ構造とみなせる ことを示し,これにより移行単位の最適化をグラフ構造におけ るノードのクラス分類問題と捉えることができた.

グラフ中のノードのクラス分類手法としてはこれまでに 数多くの手法が提案されている.まず,教師無しのグラフ ノードクラスタリング法として,ラプラシアン行列から得 られるグラフスペクトルを用いた手法や,スピングラス法 [Reichardt and Bornholdt 2006] 等の手法が知られている.教 師無しのアプローチは,教師データを要しないという点でコス ト面で優位であるものの,人が意図する移行単位とは異なる分 類がなされる可能性が高く,実用に適さない場合もある.

一方,少量の教師データを用意することで,残るノードを 人の意図に沿った形で分類できる半教師有り学習法も多く提 案されており,ラベル伝搬法やランダムウォークに基づく手法 [Talukdar et al. 2008] が知られている.本稿でもシステム開 発者の意図に沿った分類を行うという目的に則り,半教師有り のアプローチを採用することとする.

近年ではさらに、グラフに対して適用可能なニューラ ルネットワークが提案されており、グラフノードの分類 や複数のグラフにおける各グラフの分類が可能となって いる [Atwood and Towsley 2016, Kipf and Welling 2017]. Kipf らは、グラフスペクトルに基づく Graph Convolutional Network による半教師有り学習法を提案している [Kipf and Welling 2017]. 一方, Atwood らは、diffusion process を用いた Diffusion-Convolutional Neural Network (DCNN)を提案している [Atwood and Towsley 2016]. 本稿 では高精度なグラフノードの分類結果を得るため、CRUD 表 によって定義されるグラフ構造に対し DCNN を適用する.

DCNN は以下の式で定義される.

$$Z = f(W^c \odot P^* X) \tag{1}$$

$$Y = \arg\max(f(W^d \odot Z)) \tag{2}$$

ここで  $W^c$ ,  $W^d$  は  $H \times F$ のモデルパラメータ, X は  $N \times F$ の全ノードに関する素性行列, Z は diffusion-convolutional activation と呼ばれ,  $N \times H \times N$  のテンソルを表す. N は全 ノード数, F は全素性数, H はホップ数をそれぞれ表す. P\* は 隣接行列を P としたときの冪級数であり,  $P^*_{ijl} = P^1_{il}$ であ らわされる. i は元ノード, l は行き先ノード, j はホップ数を 表す. DCNN における diffusion とは, あるノードの素性値 がホップ数と冪級数に基づいて拡散することを指し, これによ り周辺ノードの影響を加味した上でグラフ中の各ノードの特徴 を捉えることが可能となる.

#### 3.2 グラフノードの分散表現

前節で述べた半教師有りグラフノード分類法を移行単位の 最適化に適用した場合,各ノードの属性情報に基づく元の表現 能力が十分でないために,学習が進まないケースが実験的に多 く見られた.本節ではグラフ埋め込み法に基づく各ノードの分 散表現を用いることでこの課題を解決する.具体的なグラフ埋 め込み法としては node2vec [Grover and Leskovec 2016] と Poincaré embedding [Nickel and Kiela 2017] を導入する. 本稿では、これら手法で得られた各ノードの分散表現を新たな 追加の属性情報として扱う. node2vec は、自然言語処理の分野で用いられる word2vec [Mikolov et al. 2013] をグラフ構造へ適用するよう拡張した手 法で,word2vec の推定に用いられる Negative Sampling を用 いた skip-gram 法の枠組みを援用し、任意のノードの近傍を 得ることで自ノードに対する分散表現を得る手法である.一列 のシーケンス構造で構成されるテキスト中の近傍とは異なり、 グラフにおける近傍の定義は様々考えられるが、node2vec で はランダムウォークによって近傍を得る.

一方 Poincaré embedding では双対空間の一種である Poincaré 空間においてグラフ埋め込みを行う手法で,ユー クリッド空間に比べ,空間を効率的に利用することができ る.特に階層構造における表現能力の高さが知られてお り,wordnet 等への適用において高い精度が報告されている [Nickel and Kiela 2017].本稿で扱うシステム移行単位の解 析においても,2.2節で導入したファイル配置情報がまさに階 層構造を取ることから,Poincaré embedding の効果が期待さ れる.

## 4. 実験

#### 4.1 実験条件

本節では提案手法の効果を確認するため,実際の段階的移行 に基づくシステム開発において,人手で設計された 8 クラス の移行単位を教師信号,および正解とみなすことで実験を行っ た.機能ノード数は 5108,テーブルノード数は 6953, CRUD 関係によるエッジは約 86000,ファイル配置情報によるエッジ は約 15000 である.CRUD 関係の種類は,通常の C,R,U,D に加え"Truncate"(T)という操作関係を含め計 5 種を扱う.

各ノードには事前に定義された 32 種類の属性の有無をそれ ぞれ付与した.人手で設計した移行単位クラスに含まれるノー ド数は,それぞれ 1106, 1542, 743, 2019, 2806, 800, 481, 2564 で,ノード数最大のクラスを常に推定結果とした場合の正解 率,チャンスレートは,21.3%となる.

また,初期に与える教師データ (シード) は,0.1%,1%,10%,30%,50%,70%,90% とデータ量を変化させつつ,全 機能・テーブルノードからランダムに選択した.さらに DCNN においては,教師データのうち,1割を validation のために使 用した.

本実験で取り扱う比較対象として,以下3つの条件の組み 合わせを行う.

- グラフ構造: CRUD 関係を単独で用いる場合 ("CRUD") と、CRUD 表にファイル配置情報を重畳させた場合 ("CRUD+filepath")の2種類を比較する.
- 半教師有りノード分類手法:半教師有りグラフノード分類 法の一種であるラベル伝搬法(LP)と, DCNNの2種類を 比較する. DCNNの epoch 数は 100 とし, validationの 誤差が最小となった時点のモデルをテストに用いた. hop 数は3, dcnn layer および dense layer 数はそれぞれ 5 層とした.
- ノードの分散表現:グラフノード分散表現を用いない ベースラインに加え, node2vec, Poincaré embedding を 128,256,512次元で試したうち,LPではテスト時最良と なった node2vecの128次元をそのまま用いることとし, DCNNでは validationのエラー率が最小となった各手法 の512次元を用いた。

評価は,各手法が推定したクラスが,人手で作成した移行単位と一致するか否かの正解率に基づいて評価した.



図 2: 推定した移行単位と人手で設計した移行単位との一致率. 横軸は学習量割合を示す. グラフ中の数値は各学習量における 最高値を表している.



図 3: CRUD 関係を持つ2つのノードが同じ移行単位クラス に属する割合を学習量毎に比較した結果. "man"は人手で設 計した移行単位の結果を示す.

#### 4.2 実験結果

各比較条件における精度を教師データ量毎に比較した結果 を図2に示す.ただし、グラフノードの分散表現を追加しな い場合については、CRUD表とファイル配置情報を併用した 場合("CRUD+filepath")でも、教師データが0.1%の時に約 23%、教師データが90%の時に約51%と、分散表現を追加 する場合と比べ精度が大きく劣ったため、見やすさのため図中 から除外している.Poincaré embeddingと node2vecの比較 では、DCNNの全ての条件下において、Poincaré embedding の方が高い精度を示した.階層構造を持つファイル配置情報 を用いる場合だけでなく、通常のCRUD表に基づく場合でも Poincaré embeddingの効果を得られた.

DCNN と LP を比較した場合,教師データ量の少ない条件 下 (0.1%~30%) で比較的高い精度が得られた.一方 LP では, 教師データ量が多い条件下 (50% 以上) において DCNN を用 いた場合よりも高い精度を示した.実利用の面において,本提 案法が特に効果を発揮するのは,教師データ量が少ない初期導 入での利用と想定され,初期に人手で用意可能な教師データの 数を 10~100 個程度と考えると,初期導入においては DCNN を用いた手法が優位と言える.その後,人手のクリーニング 等により移行単位の整理を進めた後,再度 LP で分類,チェッ クを行うことで,より実用的な効果を得られることが示唆さ れる.

次にファイル配置情報の効果について, "CRUD"と "CRUD+filepath"の結果を比較すると,ほとんどのケース においてファイル配置情報を併用した場合に精度向上の効果を もたらしている.しかし,教師データ量 0.1%の DCNN にお いてのみ,ファイル配置情報を用いない方が精度が高い傾向が 見られた.0.1%は10個のデータしかなく,教師データの選 択方法に精度が敏感に反応した可能性が考えられるが,詳細な 検証は今後の課題である

最後に提案手法によって得られた移行単位分類結果をもと に、CRUD 関係にのみ着目した最適性を評価する.図3は、各 CRUD 関係 ("C", "R", "U", "D", "T") で結ばれた機能とテー ブルの対が,同じ移行単位クラスに属する割合を示す.まず, 人手で作成した正解となる移行単位の場合 ("man"),機能と テーブルのクラス一致割合は全体平均で47.8%であった.人 は CRUD 関係以外の要因も加味した上で移行単位を決定する ため、CRUD 関係における最適性を満たさなくなったものと 考えられる.これに比べ提案手法では、教師データ量が少な い場合,比較的高い一致率を示し,特に学習量が0.1%におい てはほぼ 100% 一致している. そして教師データ量が増える に従い、人手で作成した場合の一致率に近づく結果 (学習量が 90% において,全体平均で 51.0% の一致率) となった.この ことから、提案手法は CRUD 関係における最適化を指向しつ つ,人が設計した移行単位に近い分類となるように学習が進ん でいると言える.

# 5. 結論

本稿では、大規模なシステム開発において必要となる移行 単位の最適化を目的に、システムを構成する機能とテーブルの 関係を表す CRUD 表に着目し、ソースコードのディレクトリ 構造や機能・テーブルの属性情報を加味した一元的なグラフと みなした上で、半教師有りグラフノード分類法、特に半教師有 りグラフニューラルネットワークの一種である DCNN を用い ることで、比較的少量の教師データ (0.1%) でも高精度な移行 単位の推定を可能とした.

今後の課題としては以下の三点が挙げられる.第一に,実際 の移行単位設計時には人が考慮しているが,現在まだ扱えてい ない情報を,推定器に組み入れることが挙げられる.例えば, 対象システムを用いた実際の業務がどのような手順で行われて いるかを知ることで,よりリスクの低い,人の分類により近い 移行単位の設計が可能となると考えられる.第二に,本稿では 既知として与えた移行単位の数を推定することで,これを解決 するためには教師無しグラフクラスタリングとの併用アプロー チが考えられる.第三に,本稿で各ノードの情報を分散表現か ら得ている箇所について,本来的にはグラフニューラルネット ワークの中で一括して得られることが望ましく,グラフニュー ラルネットワーク自体及び適用法の改良が必要となる.これら の課題を解決することで,より高精度かつ利便性の高い移行単 位の最適化が可能となり,ひいては情報システムの頑健性を一 層高められると考える.

# 参考文献

[木原他 2015] 木原 純平, 鈴木 邦彦, 小倉 孝昭, 浜崎 俊行, レ ガシーマイグレーションプロジェクトにおける提供品質 確保の取り組み, 2015, プロジェクトマネジメント学会誌 17(5), pages 9-14.

- [Atwood and Towsley 2016] James Atwood and Don Towsley, Diffusion-Convolutional Neural Networks, 2016, In Proceedings of the Advances in Neural Information Processing Systems 29 (NIPS 2016), pages1993–2001.
- [Kipf and Welling 2017] Thomas N. Kipf and Max Welling, Semi-Supervised Classification With Graph Convolutional Networks, 2017, In Proceedings of the 5th International Conference on Learning Representations (ICLR 2017).
- [Grover and Leskovec 2016] Aditya Grover and Jure Leskovec, node2vec: Scalable Feature Learning for Networks, 2016, In *Proceedings of the 22nd ACM* SIGKDD (KDD 2016), pages 855-864.
- [Reichardt and Bornholdt 2006] Jörg Reichardt and Stefan Bornholdt, Statistical mechanics of community detection, 2006, *Physical Review*, 74. 1
- [Talukdar et al. 2008] Partha P. Talukdar, Joseph Reisinger, Marius Paşca, Deepak Ravichandran, Rahul Bhagat and Fernando Pereira, Weaklysupervised acquisition of labeled class instances using graph random walks, 2008, In Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP 2008), pages 582-590.
- [Nickel and Kiela 2017] Maximilian Nickel and Douwe Kiela, Poincaré Embeddings for Learning Hierarchical Representations, 2017, In Proceedings of the Advances in Neural Information Processing Systems 30 (NIPS 2017), pages 6338–6347.
- [Mikolov et al. 2013] Tomas Mikolov, Kai Chen, Greg Corrado and Jeffrey Dean, Efficient estimation of word representations in vector space, 2013, arXiv preprint, arXiv:1301.3781.