

Wildcard 許容頻出部分グラフパターンのグラフ分類への応用

Applications of Frequent Subgraphs with Wildcards to Graph Classification

岡崎 文哉^{*1} 瀧川 一学^{*2*3}

Fumiya Okazaki Ichigaku Takigawa

^{*1}北海道大学工学部

School of Engineering, Hokkaido University

^{*2}北海道大学大学院情報科学研究科

Graduate School of Information Science and Technology, Hokkaido University

^{*3}科学技術振興機構, さきがけ

JST,PRESTO

In this paper, we empirically study the use of frequent subgraphs with wildcards for graph classification. Incorporating wildcards into subgraph feature representations enables to provide more flexible indicators by smaller subgraph patterns, which are obtained in a shallower search. Our experimental evaluations show the improvement on the classification accuracy for the majority of 13 used datasets. We also investigate and discuss the parameter dependency and comparisons to subgraph patterns without wildcards.

1. はじめに

グラフは, 低分子化合物の構造式 [Takigawa 13], RNA 二次構造 [Karklin 05], 自然言語処理 [Kudo 05] などの知識処理に幅広く用いられている重要なデータ構造である. 近年こうした科学分野のグラフデータが PubChem など公共のデータベースに蓄積・整備されるようになり, 有効な利活用が喫緊の課題となっている. 特に, グラフデータからの教師付き学習は, 生命科学や物質科学における構造活性相関や構造物性相関の定量的モデルとして機械学習分野で研究されており, より高精度で高効率な手法が求められている [瀧川 16].

本稿では, グラフ分類として, グラフデータからの教師付き学習を扱う. グラフデータからの教師付き学習では, 特徴量として部分グラフの有無を用いることが多く, この部分グラフの特徴選択が学習モデルの精度を左右する. しかし, 与えられたグラフデータに起こるすべての部分グラフを列挙することは現実的に不可能である. したがって, グラフカーネル法 [Kashima 03] や ECFP 法 [Rogers 10] など, 対象となる部分グラフのクラスを予め発見的に制限する手法や, gBoost 法 [Saigo 09] や Adaboost に基づく手法 [Kudo 05] など, 必要な部分グラフのみを効率的に探索しながら学習を行う手法が提案されている. 全ての部分グラフの列挙は現実的に不可能であるため, 特徴として使う頻出部分グラフは, 閾値 σ (与えられたグラフデータベースに対して σ 以上のグラフに出現する) を与え, 列挙することを考える. その時, σ を小さくしたほうが特徴候補の数を増やすことができるため精度向上が期待できる. しかし, σ を下げると出力数が指数的に増加するため, 深い探索が必要な事例では有効な特徴がある場合でもメモリや時間に対するトレードオフがある.

本稿では, 我々が提案した wildcard を許容した頻出部分グラフマイニング [岡崎 15] のグラフ分類への応用を提案する. グラフ分類において使用する部分グラフ特徴は, 包含関係に由来する冗長性や相関を持っており, wildcard を許容した緩和特徴表現により, 厳密表現と比べて小さい部分グラフで柔軟な記述力を得ることが期待できる. そして, wildcard 許容頻出

部分グラフを特徴に用いた場合の効果を確認するため, 13 種類の低分子化合物の構造活性相関データを用いて, 実験的に検証する. また, 通常の頻出部分グラフを特徴に用いた場合やパラメータを変化させた場合と比較し, wildcard 許容頻出部分グラフの有効性と条件の分析を行う.

第 2 節では, 頻出部分グラフを用いたグラフ分類, 第 3 節では, グラフ分類に用いた機械学習の手法 Random Forest の説明を行う. 第 4 節では, wildcard 許容頻出部分グラフによるグラフ分類を提案する. 第 5 節では, 提案手法に対する実験を示し, 第 6 節で, 実験に対する考察を行う. 第 7 節で, 結論を述べる.

2. 頻出部分グラフによるグラフ分類

本稿において, グラフ分類とは, グラフ集合に対する教師付き学習を指す. 複数の種類のグラフ集合が与えられた時, これらのグラフ集合を分けるための特徴付けや与えられたグラフがいずれの集合に属するかを判別する問題である.

本稿では, 2 種類のグラフ集合を入力として与え, 分類器を構築する. 分類器の構築において, 特徴量の決定が重要な課題である. この特徴量の決定に対して部分パス, 部分木, 部分グラフなど様々な特徴量を考えることができるが, 我々は頻出部分グラフによる特徴付けに着目した.

すべての部分グラフパターンの列挙は, 出力数が組み合わせ爆発を起こすため, 現実的に求めることは困難である. そのため, 通常, ある閾値を与え, σ 個以上のグラフに出現するパターンを頻出部分グラフとして列挙する. この閾値 σ を最小支持度 (以下, minsup) と呼ぶ.

頻出部分グラフによるグラフ分類では, minsup によって精度が異なる. 原理上は, minsup を下げることで, 得られる部分グラフパターンの集合は単調に増大するので, 用いる機械学習アルゴリズムが適切に対応できれば, 精度も単調に向上する. しかし, minsup を下げれば出力数 (用いる特徴の数) は指数的に増加するため, 計算時間や使用メモリの点で現実的には全列挙が困難となる. そのため, 計算可能な特徴量の数には制限があり, この制限に対して分類の精度を上げる必要がある. このような場合, 従来は用いるグラフの型をパスや木のみで制限する方法が検討されてきた. 本稿では型はグラフのまま変え

ず, wildcard を許容した特徴表現を用いた場合の有効性について実験的に検証を行う。

3. Random Forest

本稿では、グラフ分類で用いる特徴集合の良さの評価が目的であるため、分類に用いる学習アルゴリズムには統一して Random Forest を用いた。

Random Forest[Breiman 01] とは、複数の決定木を統合してマルチクラス分類を行うアンサンブル学習アルゴリズムである。このアルゴリズムの特徴は、ブートストラップを取り入れ過学習を防ぐことが可能な点、Random Feature Selection[Ho 98] を取り入れることで特徴ベクトルの次元数が大きくても高速に学習が可能である点にある。これらのメリットを持ち、クラス識別、クラスタリング、回帰分析を行うことができる。このように高い汎用性を持つため、様々な分野に応用されている。

本稿において、Random Forest を用いた理由は、以上の Random Forest の特性に加えて、各特徴に対して重要度を定めることができ、どの特徴が分類に寄与していたかの説明が定量的にできるためである。

4. 提案手法

本稿では、グラフ分類のための特徴量として、wildcard 許容頻出部分グラフ [岡崎 15] を用いる手法を提案する。本節では、wildcard 許容頻出部分グラフの列挙対象について説明する。

定義 1 (wildcard 許容頻出部分グラフ) wildcard を k 個許容した頻出部分グラフは、任意のラベルにマッチするラベルである wildcard を k 個以下含む部分グラフで、入力で与えた閾値 $minsup$ 以上のグラフに含まれる部分グラフ全てである。

一般的なグラフデータにおいては辺のラベルにおける wildcard が意味を持つことがあると考えられるが、本稿において、取り扱うデータベースは、化合物の分子構造を表したものであり、グラフの辺のラベルは頂点間の結合数を表しているため、この部分に対する wildcard 許容は重要な構造になる場合が少なくないと判断することができる。そのため、本稿の列挙対象、その対象実験においては、頂点ラベルに対してのみ wildcard 許容を考える。

頻出部分グラフは $minsup$ を下げることで、列挙対象となるパターン数は膨大に増加する。さらに、wildcard を許容した場合、許容数を増やすことで大幅に列挙数が増加し、適切な対処をしなければ、現実的に列挙できない数になる。我々は、wildcard を許容した頻出部分グラフマイニング [岡崎 15] において、効率的な列挙の方法を提案している。さらに、出力数の増加問題に対して、冗長なパターンとその削減方法、その頻出飽和・極大パターン集合を求める手法を示した。

本稿では、wildcard 許容頻出部分グラフをグラフ分類に対する特徴として使う手法を提案し、wildcard の有効性について検証する。また、単純に wildcard 許容頻出部分グラフを用いる手法に加え、 $minsup$ が大きい部分に wildcard を許容し、出力数の増加を抑えるために、通常の頻出部分グラフと組み合わせるといふ、wildcard 許容の応用を提案し、実験を行う。さらに、冗長なパターンの削減や頻出飽和・極大パターン集合を用いることで、応用するグラフデータベースに対して wildcard 許容頻出部分グラフをうまく利用することができる。と期待できる。

グラフ分類の特徴量として、wildcard 許容頻出部分グラフを用いる場合、どのパターンが重要なパターンであるかという

CPU	Intel(R) Core(TM) i7-3770K 3.50GHz
メモリ	33GB
OS	Ubuntu 14.04.3 LTS

表 1: 実験環境

データセット名	グラフ数	平均ノード数	平均エッジ数
CPDB	684	14.1	14.6
Mutag	188	17.9	19.8
NCI1	1000	35.3	38.5
NCI41	1000	35.5	38.7
NCI47	1000	35.5	38.7
NCI81	1000	34.4	37.5
NCI83	1000	34.5	37.5
NCI109	1000	35.2	38.4
NCI123	1000	34.2	37.3
NCI145	1000	34.9	38.0
NCI167	1000	30.5	33.5
NCI220	1000	28.3	30.2
NCI330	1000	30.5	33.0

表 2: 使用したデータセット

解釈が困難である。本稿における実験では、wildcard を許容することによる有効性に対して、一般的な知見を得るために、各データベースによる設定の変更は行わず、冗長なパターンの削減もせず比較を行った。次の節で、詳細について議論する。

5. 実験

本節では、wildcard 許容頻出部分グラフをグラフ分類の特徴として使用する実験を示す。本稿では、13 種類のデータベースに対して、通常の頻出部分グラフと wildcard 許容頻出部分グラフを比較することで、wildcard 許容部分グラフの有効性について検証を行う。

5.1 実験環境とデータセット

本実験は表 2 に示す 13 種類のグラフデータを用いて、表 1 に示す環境で行った。データセット名「NCIx」は PubChem BioAssay の AID 番号 x のデータセットである*1。これらのデータセットのうちグラフ数が膨大なものは、各データにおける 1000 個のランダムサブセットを用いた。その時、正例と負例の数が同じになるように選択を行った。その際に、正例が 500 に満たないデータに対しては、グラフ数が 1000 になるように負例を追加した。構造活性相関の実験の実データでは通常、負例数が極めて多く、この不均衡性 (imbalance) を考慮した評価尺度を用いる必要がある。例えば、ROC50[Wale 08] を用いるなどが提案されているものの、上位 50 個で良いかなど解釈の問題のため、広く確立した尺度とは言い難い。本稿では単純に正例と負例の数を均衡化し、通常の高正解率により特徴集合の良さを評価する。

5.2 実験結果

本実験において、正解率はデータを 10-分割交差検証を行った際のテストデータに対する正解率 (ACC と表記) を示す。用いたデータセット 13 種類のうち、CPDB, Mutag は $minsup$ を変化させた時の精度の変動を示すために使用した。また、NCIx

*1 <http://www.ncbi.nlm.nih.gov/pcassay>

のデータを用いて、wildcard を許容した場合の精度と出力数に対する傾向を検証するための実験を行った。

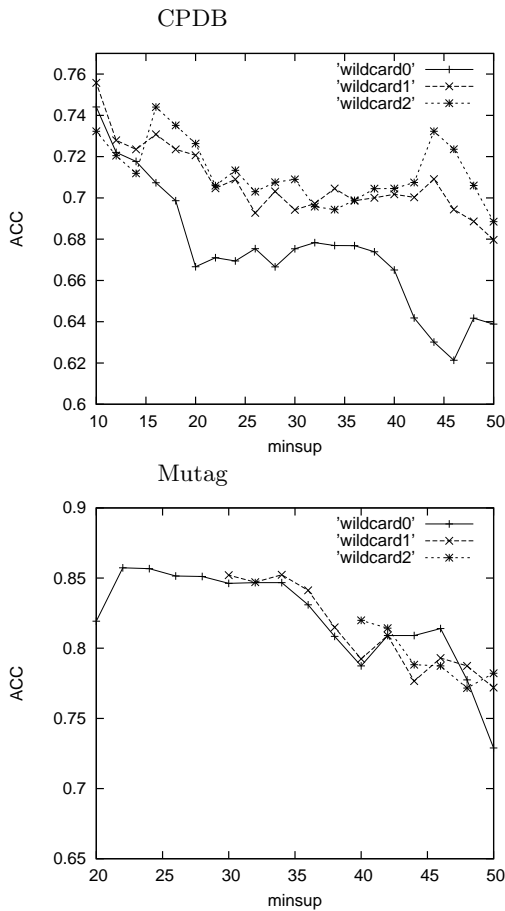


図 1: minsup の変化と正解率の実験結果

図 1 に wildcard 許容に対する minsup を変動させた時の実験結果を示す。上段が CPDB, 下段が Mutag に対する実験結果である。この図において、「wildcardx」は wildcard を x 個許容を指す。

図 1 上段の CPDB では、minsup が全体的に精度の向上が見られる。minsup を十分に下げるとあまり精度に差がないが、通常の頻出部分グラフパターンによる結果と同等か、それ以上の結果が出ている。また、wildcard0 (wildcard 許容なし) では、minsup が 20%以上になると 0.72-0.74 から 0.66-0.68 に急に精度が下がっている。しかし、wildcard を許容することで minsup が 20%以上でも 0.70-0.72 の精度が得られている。

図 1 下段の Mutag では、minsup が 20%-33%では精度にほぼ変化がなく、33 以上になると、wildcard の許容にかかわらず、0.75-0.8 周辺に精度が落ち込む。wildcard を許容することによる精度の悪化は見られないが、精度の向上も見られないことから、wildcard の有効なデータベースではないことがわかる。

表 3 に、NCI データセットに対する実験結果を示す。この表では、各データセットに対して、通常の頻出部分グラフ、wildcard を 1 つ許容した頻出部分グラフ、wildcard を 2 つ許容した頻出部分グラフについて、minsup を 45%から 5%まで 5%ずつ変化させた場合の出力数と正解率 (ACC と表記) を示す。NA と表記した部分は現実環境でメモリ不足により求めることができなかったものである。このとき、精度が一番良かったも

のを太字としている。このデータセットにおいて、11 個中 10 個で wildcard 許容のほうが良い結果が出ている。

さらに、右側には、wildcard 許容と通常の頻出部分グラフを組み合わせた場合の実験結果を示している。これは先ほど太字で示した部分の特徴に加えて、それより minsup が小さい頻出部分グラフを追加した場合の出力数と正解率である。例えば、NCI1 では、wildcard2 許容の 30%で良い結果が出ている。この時の特徴量 10721 に加えて、通常の頻出部分グラフを minsup10%まで加える。ここで加えている特徴量は 30% 10%における (3283-225) 個である。それにより、特徴量は 13778 個であり、精度は wildcard2 許容の 30%に対して向上しているという結果を得た。このように、組み合わせにより、精度が向上したのに対して太字で示している。実験結果より、wildcard 許容により精度が向上した 10 個のデータセットのうち、3 個のデータセットで、組み合わせによる精度の向上が見られた。

これらの結果に加え、Random Forest による特徴量の重要度の高かったものを見てみたところ、wildcard 許容部分グラフが上位の 20 個程度をほぼ占めていた。これにより、通常の頻出部分グラフより、wildcard 許容部分グラフのほうが分類に寄与する特徴として Random Forest で選ばれやすかったことが分かる。従って、wildcard 許容による拡張で得られた部分グラフ特徴は冗長なものではなく実際に分類に対する有効性を持つものと言える。

6. 考察

実験結果より、wildcard を許容することで、単純に minsup を下げるより、精度が向上するデータベースが多いという結果になった。一般的に、Random Forest は弱学習器として決定木をもつため、正例と負例をよりうまく分ける特徴が多く存在するほど精度が上がる場合が多い。wildcard を許容したパターンは通常のパターンに加え、データベースをうまく表現するような特徴としての候補を得ることができ、このようなパターンが精度の向上に寄与すると考えられる。

また、頻出部分グラフでは、minsup を大きく設定すると、得られるパターン数が少なく表現力が弱い。それに対し、wildcard を許容すると、minsup を大きく設定した場合に表現力が厚みを増すため、精度が向上すると考えられる。しかし、wildcard を許容することによる出力数の増加が原因で minsup を下げて計算することが困難である。そこで、minsup を下げた部分の表現を与えるために頻出部分グラフと組み合わせることで、さらに精度が向上するデータベースが存在することも実験により示すことができた。

7. 結論と今後の課題

本稿では、グラフ分類に対する特徴量として、wildcard 許容頻出部分グラフを利用する手法を提案し、実験により wildcard 許容の有効性について検証し、考察を行った。その結果、wildcard が有効であるデータベースとそうでないデータベースが存在するが、wildcard を許容した場合に精度の向上が見られることが多いという実験結果となった。これは、wildcard を許容することによる表現力の向上の結果であると言える。実際に、応用する際には、最小支持度の設定、冗長パターンの削除や頻出飽和・極大パターン集合の利用、wildcard 許容と許容なしの組み合わせなど、適切な設定を行うことで、精度の向上に務める必要がある。適切な設定や wildcard が有効なデータベースの性質などに対する研究が今後の課題である。

