

ビッグデータ分析技術ワークショップ
～大規模グラフマイニング技術と応用～



高速グラフマイニング技術 Grapon による大規模グラフ処理と応用

NTTソフトウェアイノベーションセンター
2017年3月5日



高速グラフマイニング技術Graponの中から今回、グラフベースの機械学習として有望なラベル伝播を中心に、適用と大規模化対応の研究をご紹介します。

<トピックス>

- はじめに
- データどうしの類似度によるグラフ構築
- ラベル伝播による文書分類
- 大規模化対応





はじめに

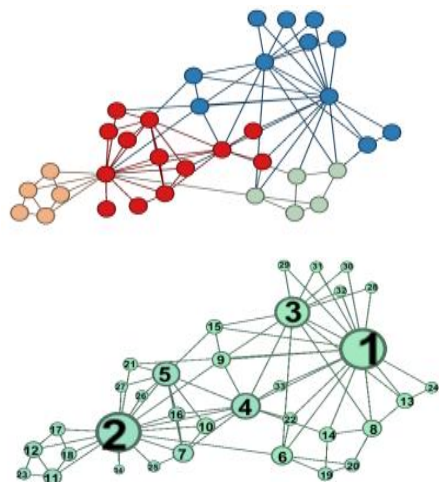
- **グラフマイニングの利点**
- **研究開発の注力分野**
- **グラフマイニング研究に関わる諸活動**

グラフマイニングの利点

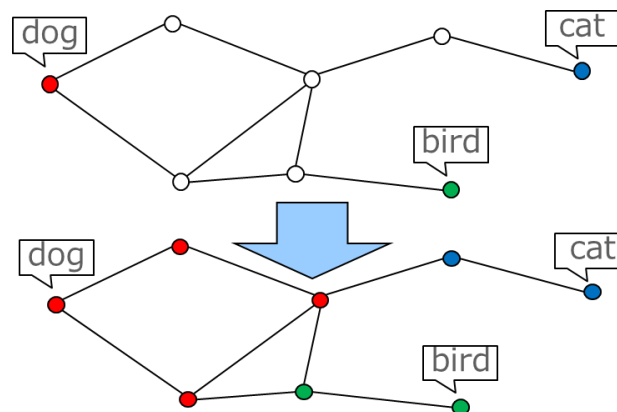


- つながり構造を**効率的に処理可能**
- グラフ構造があることで**教師データが少なくても高精度な分析が可能**
- **可視化**との親和性が高い

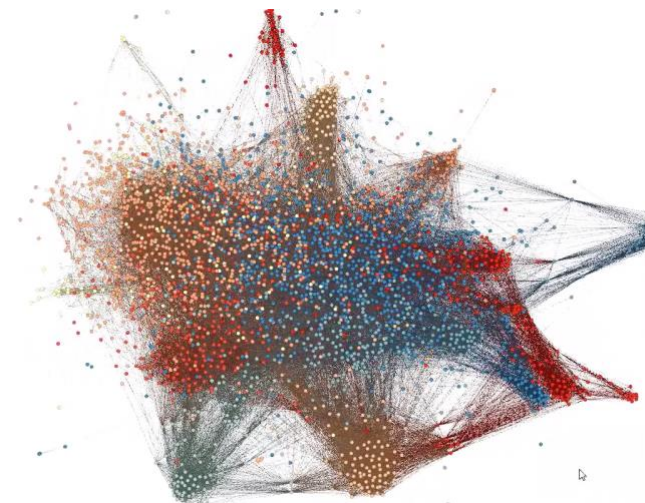
つながりを効率的に処理



少ない教師データから高精度に分析
(ラベル伝播)



可視化による情報把握の容易化
特徴的な部分をすばやく発見



グラフを活用した分析は有望だが、大規模グラフの分析には膨大な時間がかかる。
そこで、活用が見込まれる分析アルゴリズムについて高速化技術を研究開発。

-> Grapon(高速な分析アルゴリズム群)

- クラスタリング, ランキング, 推定・分類, グラフ構築, 効率的な並列化など
- アルゴリズムの改良で数10倍から数100倍の高速化
- 1億ノード規模以上のグラフのスケールラブルな処理を目指して取組中。

<主な研究成果>

- 高速Modularityクラスタリング AAAI'13
- 高速PageRank AAAI'13
- 高速Personalized PageRank SIGMOD'13
- 高速Label Propagation ICML'14
- 高速グラフ分割（等粒度クラスタリング）DEIM'15
- 高速SCAN:Structural Clustering Algorithm for Network PVLDB(2015)
- 高速Belief Propagation IJCAI'15
- 高速な並列グラフ処理（Rabbit Order）IPDPS'16
- 高速L1 sparse graph構築 PVLDB(2016)

グラフマイニング研究に関する諸活動



価値ある技術の創出と、コラボレーションを通じた検証



データどうしの類似度によるグラフ構築

- グラフを活用した分析例
- グラフ構築手法
- 構築したグラフの概観例

グラフを活用した分析例



グラフ構造が自明でない場合は，データどうしの類似度からグラフを構築

グラフ構造が自明な場合，その構造を利用（交友関係，道路の接続関係など）

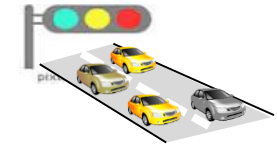
SNSにおけるインフルエンサー発見



交通管制に向けた広域道路網の最適メッシュ分割



予測・信号制御・渋滞回避

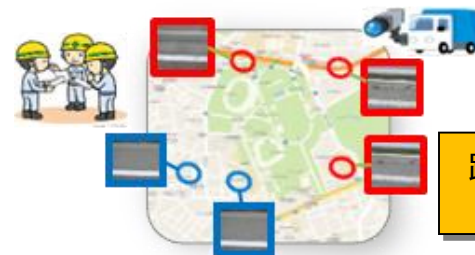


グラフ構造が自明でない場合，多次元量から良いグラフを作る必要がある（センサーデータ，画像，文書など）

状態どうしの類似性を計算し，カテゴリに分類



画像どうしの類似性から，画像をカテゴリに分類



グラフ構築にはk-NNのほか、**全体のバランスを考慮したグラフを構築するものとして、Linear Neighborhood, L1 sparse**等が提案されている。

- k-NNグラフ

近傍の上位 k ノードどうしでエッジを張る。

- Linear Neighborhoodグラフ

あるノード \vec{x}_i の近傍ノードとの線形和 $\sum W_{ij} \vec{x}_j$ を計算。これとの距離が小さくなるように重み W_{ij} を決める。

$$\min \varepsilon = \sum_i \left\| \underset{\substack{\uparrow \\ D \times 1}}{\vec{x}_i} - \sum_j \underset{\substack{\uparrow \\ \text{重み(スカラー値)}}}{W_{ij}} \underset{\substack{\uparrow \\ D \times 1}}{\vec{x}_j} \right\|^2$$

凡例

N : ノード数

D : 次元数

- L1 sparseグラフ

最適化問題を解くことでエッジ有無を自動的に決定。ノイズに強く、データ依存性が少ない。

ただしグラフのノード数よりも（ノードの）次元数が非常に多い場合に有効。

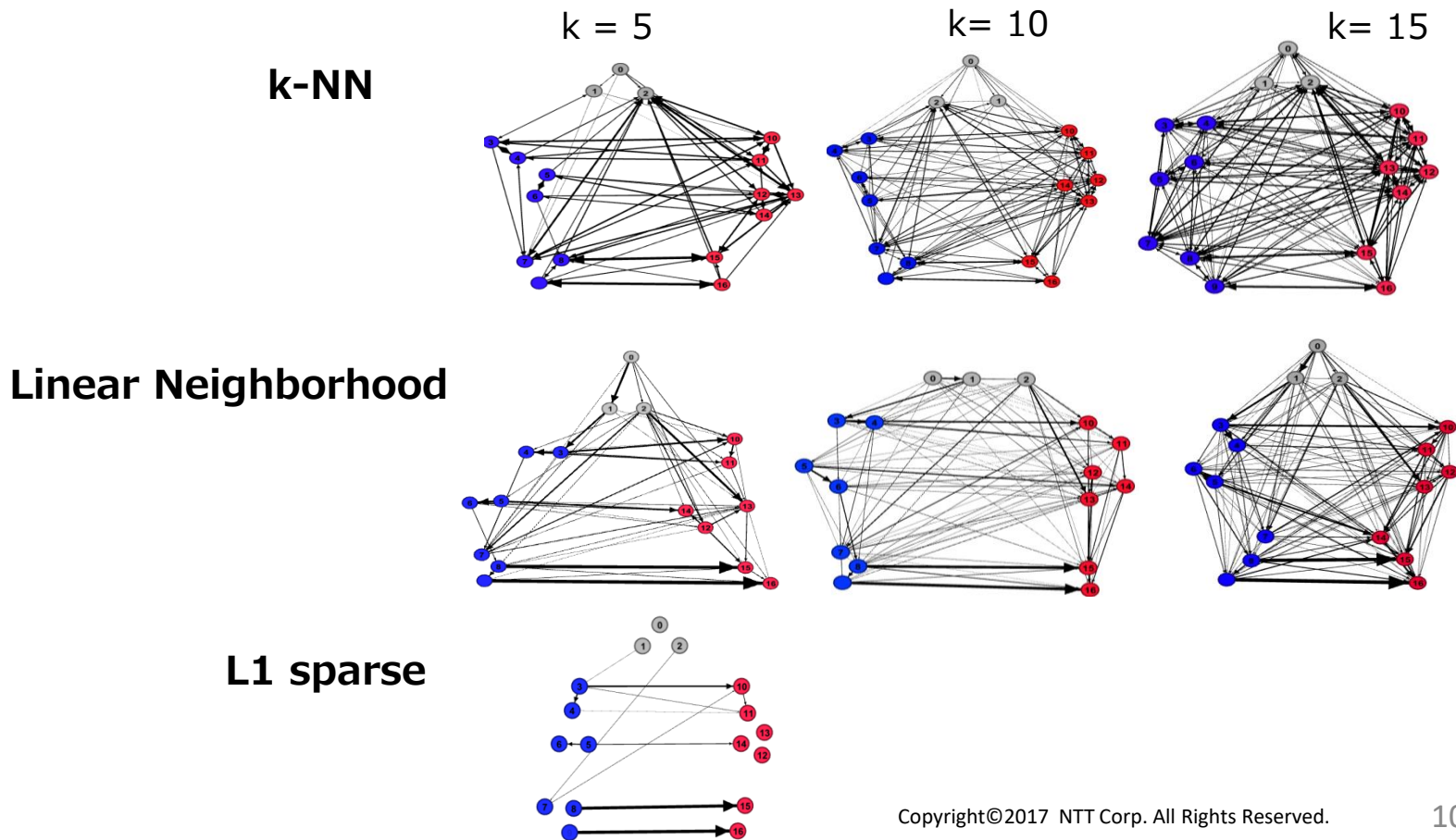
$$\min_{\vec{w}_p \in \mathcal{R}^N} \frac{1}{2M} \left\| \underset{\substack{\uparrow \\ 1 \times D}}{\vec{X}_p} - \underset{\substack{\uparrow \\ \text{重みベクトル} \\ (1 \times N)}}{\vec{w}_p} \underset{\substack{\uparrow \\ N \times D}}{X} \right\|_2^2 + \lambda \left\| \underset{\substack{\uparrow \\ \text{重みベクトルのL1ノルム(ノード}p\text{と他ノードとのエッジ重みの絶対値和)}}{\vec{w}_p} \right\|_1$$

構築したグラフの概観例



Linear NeighboringやL1 sparseは疎なグラフになりやすい傾向。

文書から単語共起による類似度で構築したグラフの例



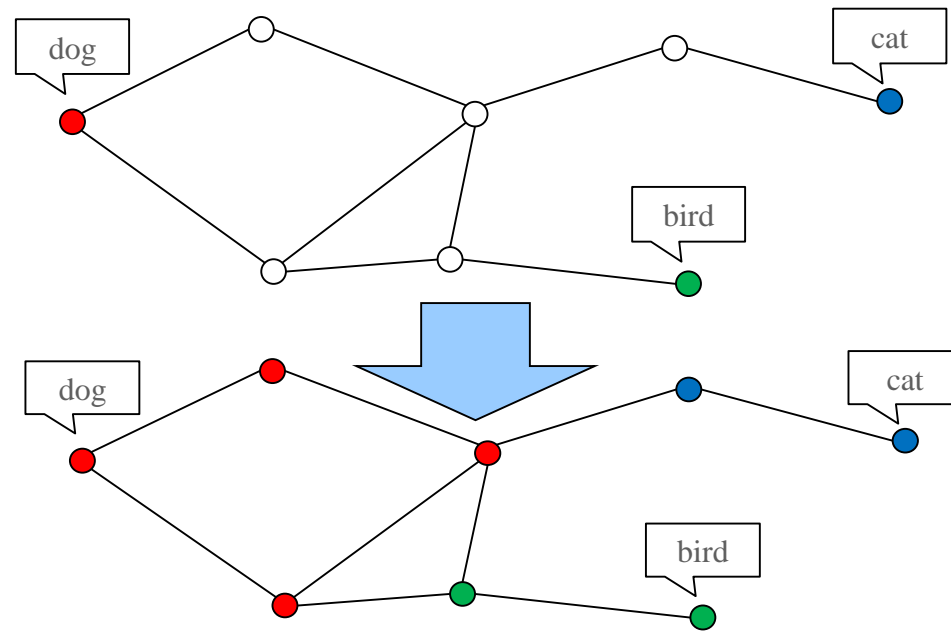


ラベル伝播による文書分類

- ラベル伝播 (Label Propagation)
- ITU勧告文書の分類



グラフベースの半教師あり学習。ラベル既知のノードから、ラベル未知のノードにラベルを推定することで、データのカテゴリを推定・分類する分析に適用可能。



$$F = (I - \alpha S)^{-1} Y$$

$n \times c$ $n \times n$ $n \times c$

n : # of nodes

c : # of category

F : score matrix

I : identity matrix

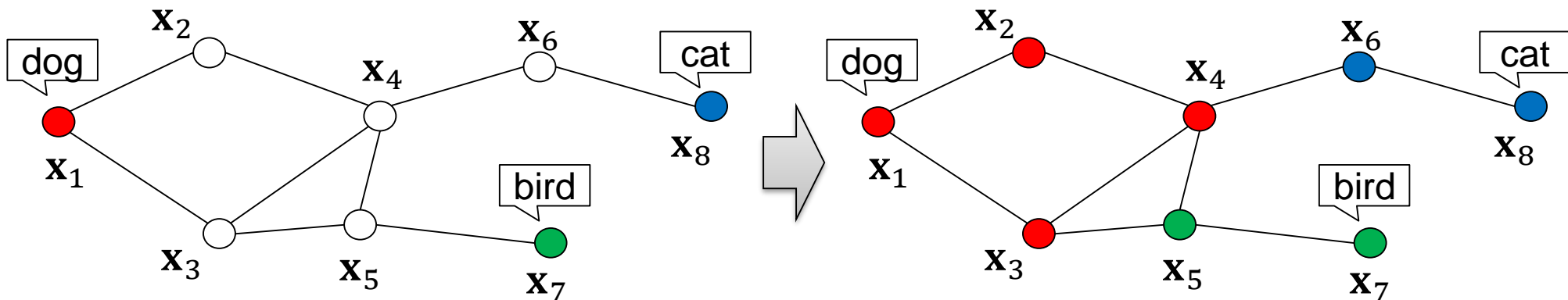
α : constant parameter

S : normalized adjacency matrix

Y : initial label matrix



各ラベルのそれぞれについて, スコアを計算していく.



Y^T

	x_1	x_2	x_3	x_4	x_5	x_6	x_7	x_8
dog	1	0	0	0	0	0	0	0
cat	0	0	0	0	0	0	0	1
bird	0	0	0	0	0	0	1	0

F^T

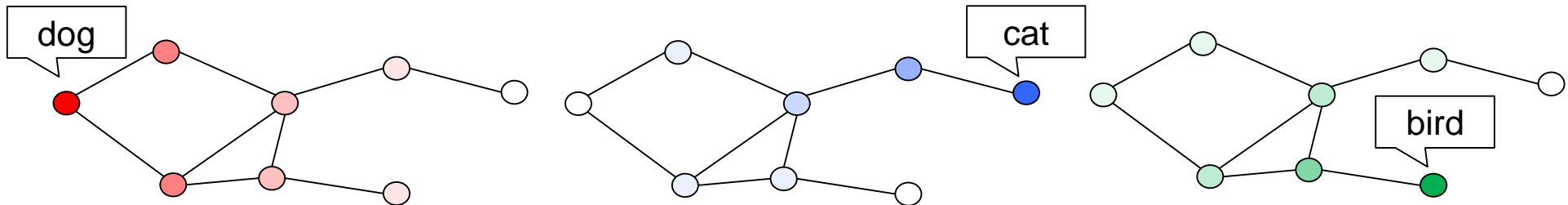
各数値は, 候補となる各ラベルのスコア

	x_1	x_2	x_3	x_4	x_5	x_6	x_7	x_8
dog	1	.9	.8	.6	.1	.1	0	0
cat	0	0	0	.3	0	.9	0	1
bird	0	.1	.2	.1	.9	0	1	0

ラベル伝播 (Label Propagation) (3/3)



- 隣接するノードのラベルどうしは、同じラベルが付与されやすい (**smoothness**)
- 教師ラベルの近傍は、教師と同じラベルが付与されやすい (**fitting**)
- α は0~1. 0に近いほど教師ラベルの影響を受けやすい, 1に近いほどグラフ構造の影響を受けやすい.



$$\min_F \frac{1}{2} \sum_{i,j} W_{ij} \|\bar{F}_{i\cdot} - \bar{F}_{j\cdot}\|^2 + \left(\frac{1}{\alpha} - 1\right) \sum_i \|F_{i\cdot} - Y_{i\cdot}\|^2$$

smoothness

fitting

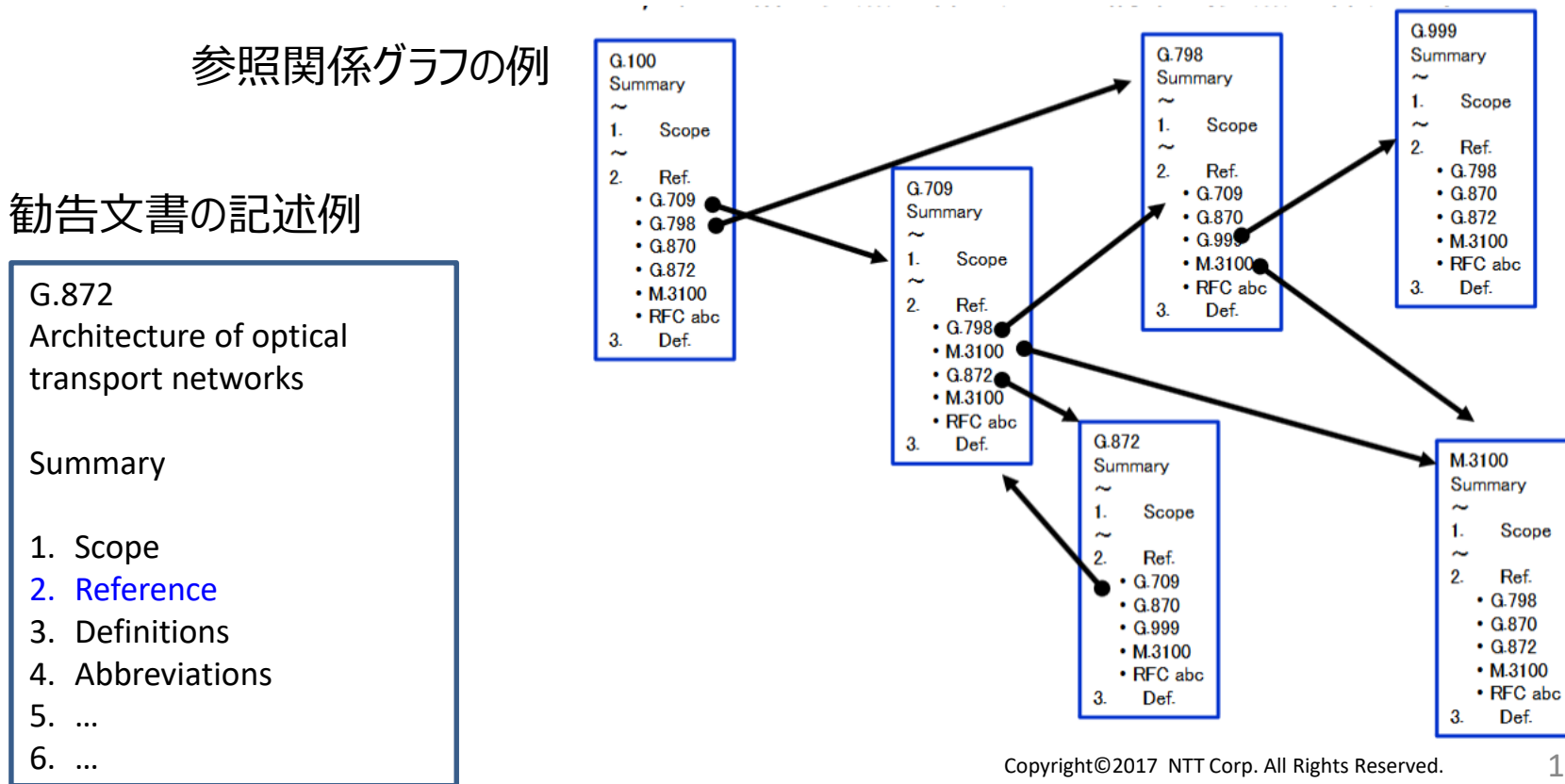
$$\Rightarrow F = (I - \alpha S)^{-1} Y, \quad S = D^{-\frac{1}{2}} W D^{-\frac{1}{2}}$$

ITU勧告の分類 (1/4)



- ITU勧告は参照関係が整備されている。 **クラスタリング**において単語共起グラフに比べて**参照関係グラフ**の精度が高いことを確認※
- 今回、**ラベル伝播による分類**でも同様かどうかを検証した

※モジュラリティ0.62で最も良い精度。他のグラフ構築手法によるモジュラリティは0.17~0.61 (昨年のビッグデータ分析技術ワークショップでご紹介)



ITU勧告の分類 (2/4)



参照関係のほか，単語共起による類似度でグラフを構築した。

- グラフ構築手法は，k-NN, Linear Neighboring, L1 sparseの3種類

ITU勧告Gシリーズに頻出する単語の例

出現数	word	
17	network	
14	requirement	synchronization
12	equipment	distribution
	packet	frequency
11	clock	time
	timing	
10	layer	
9	Ethernet	element
	methods	IEEE
	architecture	phase physical
8	case protocol	
7	performance	tolerance
	signal	details
	further	precision
6	minimum	PTP
	study	support
	types	application
6	relevant	aspects
	future	interoperability
	specification	telecom
5	layers	transport
	IP	slave
	media	consistent
	function	first
	reference devices	functions operation

出現数	word				
4	SDH	interfaces	conditions	definitions	
	information	limited	environmental	delivery	
	different	limits	generation	environment	
	order	MPLS	normal	G.8260	
	required	RFC	related	manner	
	currently	scenarios	transfer	necessary	
	G.8110	particular	applicable	parameters	
	IETF	boundary	packet-based	utilize	
	3	measurements	manufacturers	addresses	characteristics
		covered	produced	noise	configuration
jitter		quality	proper	design	
G.8261		satisfactory	timed	full	
client		exceeded	base	G.8275	
end		method	covers	messages	
adhered		accuracy error	form	profile	
2	allocation	G.803	contain	G.8265	
	compliance	NEs	level	G.8265/Y.1365	
	digital	Option	allowed	mode	
	G.709	bandwidth	International	carrier	
	G.798	connection	located	general	
	individual	delay	national	interaction	
	ODUK	G.8261/Y.1361	output	nodes	
	optical	holdover	system	packetbased	
	OTN	interworking	detailed	algorithm	
	rate	PDV	framework	best	
	respective	period	measurement	complement	
	structure	second	applies	exchange	
	Additional	variation	EEOption	G.8275/Y.1369	
	TDM	format	G.813	high-level	
	CES	technology	G.8271	mapping	
hypothetical	versions	operates	modes		
metrics	NTP	T-GM	options		

ITU勧告の分類 (3/4)



- 分析対象は、ITU-TのGシリーズ(伝送システム及びメディア、デジタルシステム及びネットワーク)の284件。参照先も含めて計750件の文書。
- これを7カテゴリに分類する。教師データは全体の10%および30%。

分析対象のITU勧告文書

カテゴリ	内容	データ数
0	G.100-199 International telephone connection and circuits	25
1	G.200-299 General characteristics common to all analogue carrier-transmission systems	1
2	G.600-699 Transmission media and optical systems characteristics	29
3	G.900-999 Digital sections and digital line system	59
4	G.1000-1999 Multimedia Quality of Service and performance – Generic and user-related aspects	13
5	G.9000-9999 Access networks	21
6	上記以外	136

※カテゴリ分けは、ITU-T標準化活動におけるラポータ（課題担当）経験者のNTTアドバンステクノロジー森田直孝氏によるもの

ITU勧告の分類 (4/4)



- ・ 参照関係のほかは，単語共起(438次元)を用いた
- ・ ラベル伝播は，教師データが少ない場合(10%)に特に他よりも**高精度**
- ・ 参照関係が**もっとも高精度**

分類手法	グラフ構築法	正解率				正解率 (平均)
		教師データ10%		教師データ30%		
SVM※1	-	57%	69%	84%	82%	73%
K-NN※2	-	43%	55%	66%	67%	58%
ラベル伝播	L1 sparse ($\alpha=0.3$)	71%	72%	88%	84%	79%
	L1 sparse ($\alpha=0.8$)	77%	74%	89%	86%	82%
	Linear Neighborhood ($\alpha=0.3$)	63%	73%	82%	83%	75%
	Linear Neighborhood ($\alpha=0.8$)	65%	77%	81%	82%	70%
	K-NN ($\alpha=0.3$)	61%	65%	77%	77%	70%
	K-NN ($\alpha=0.8$)	67%	64%	71%	70%	68%
	参照関係 ($\alpha=0.3$)	80%	81%	91%	91%	85%
	参照関係 ($\alpha=0.8$)	80%	81%	91%	90%	85%

※1 SVC with linear kernel

※2 K=5. ユークリッド距離



大規模化対応

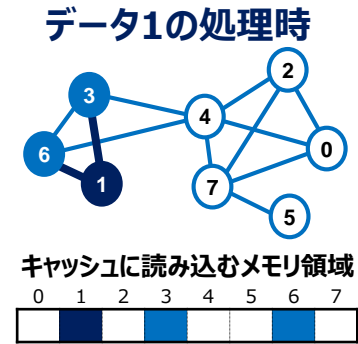
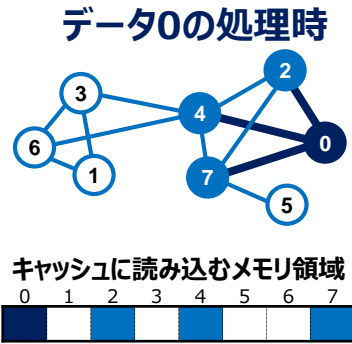
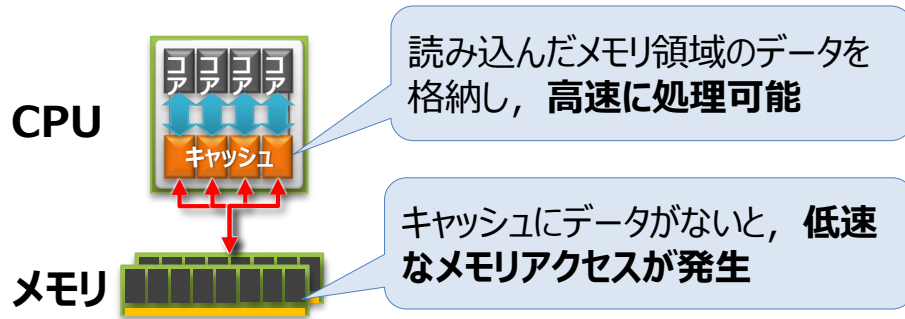
- グラフのReordering
- Reorderingによる効率的な並列処理
- 並列処理によるラベル伝播の高速化例

グラフのReordering

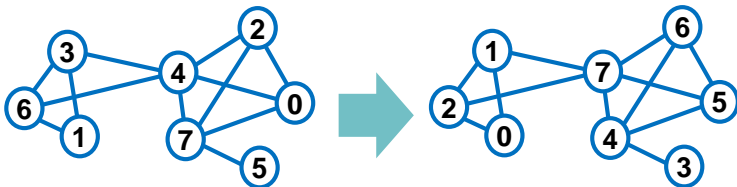
データのID番号の並びを最適化(reorder)してキャッシュヒット率を多くする。

グラフマイニングにおけるメモリアクセス

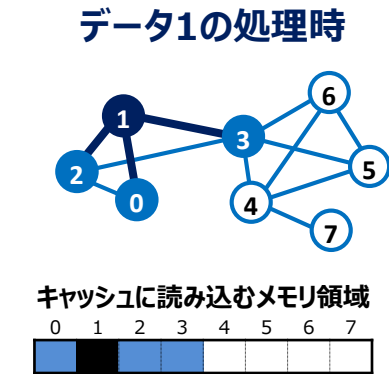
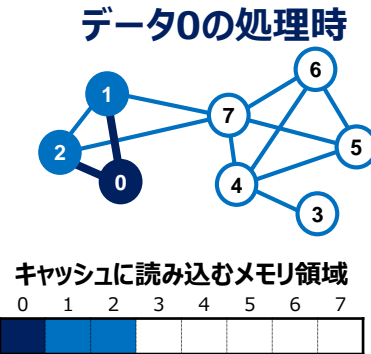
隣接するデータを順にメモリから読み込んでいく



本技術: データの並びを最適化



データのID番号を振りなおし、隣接するデータ
どうしをメモリ上の近傍に配置

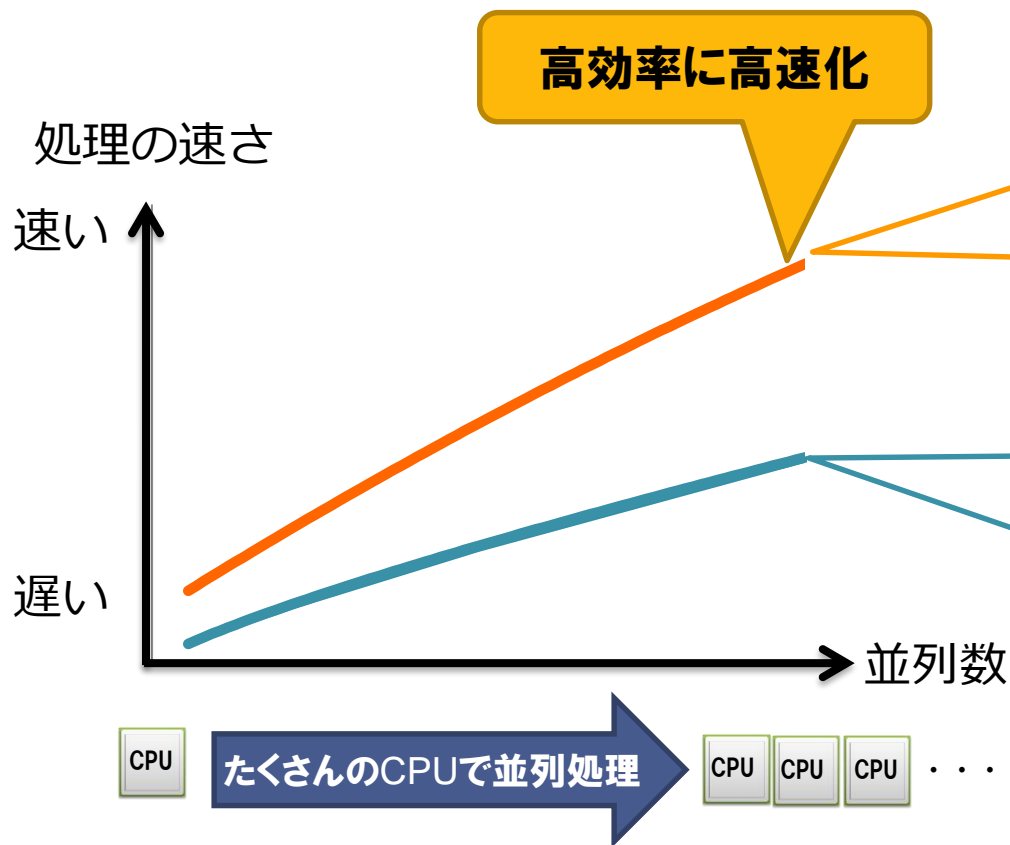


キャッシュを最大限活用
不要なメモリアクセスとCPU間通信を低減

Reorderingによる効率的な並列処理

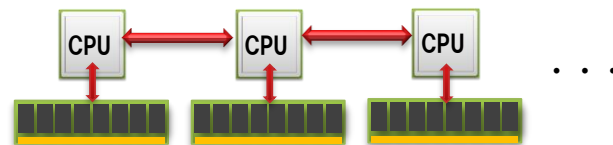


グラフ処理の並列化により、要件に合わせた高速化を可能とする。
メモリバンド幅のボトルネックを回避する高いスケール性能を目指す。



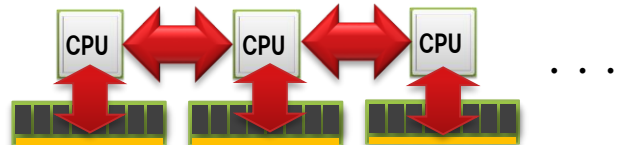
Reorderによる効率的な並列処理

メモリアクセスとCPU間通信を抑えられ、並列数にあわせたスケールが可能



グラフ処理の並列化 (reorder前)

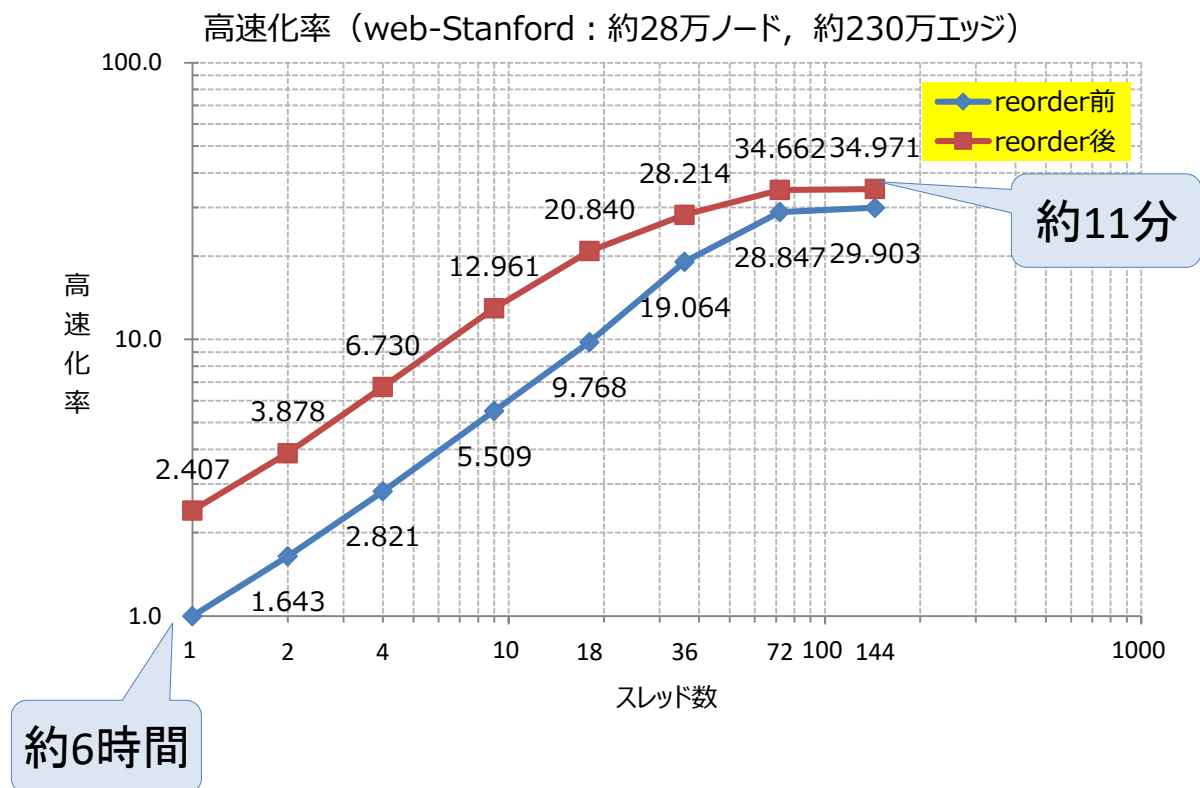
ただし、このままではメモリアクセスとCPU間通信が多く発生し、なかなかスケールしない可能性



並列処理によるラベル伝播の高速化例



- Web-Stanfordデータについて，72コア環境で35倍高速化
- 並列化にはOpenMPを利用



測定環境 : Xeon E7-8890v3 (物理18コア) x 4ソケット / メモリ2TB



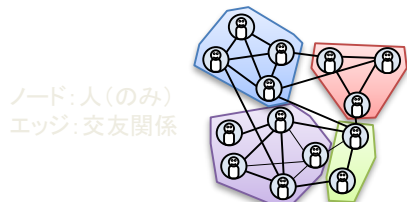
まとめと今後の予定

- **グラフマイニングの適用例として、ITU勧告文書群からグラフを構築し、ラベル伝播を行うことで、少ない教師データで高精度に分類できることを示した。**
- **グラフマイニングの大規模化対応として、効率的な並列処理を行うことでラベル伝播を高効率にスケール可能であることを示した。**

多様な属性を持つグラフデータを分析可能とし、さまざまな現実の課題解決につなげる。

ソーシャルグラフ

属性を無視した同質なノードとエッジに縮約
エッジは直感的なものを利用



さまざまなセンサからのデータ分析・処理

たとえば、データをグラフ構造として捉えることで効率的な処理ができる分野の開拓を目指す

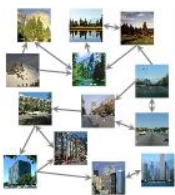


確率伝播による領域分割



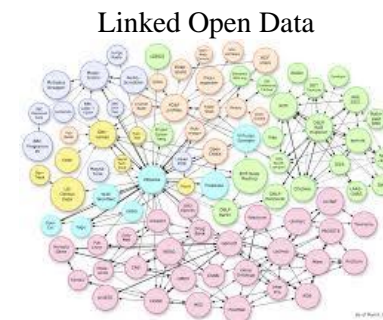
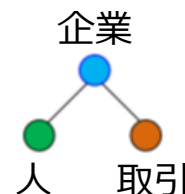
画像どうしの類似度など 多様なグラフ構造をマイニング

たとえば、グラフ構造を与えることで、教師データが少ない場合でも高精度にクラスタリング・分類するタスクを扱うことを目指す



プロパティグラフからの検索・マイニング

たとえば、属性を考慮することで、多様な類似構造の発見・検索等のタスクを扱うことを目指す





Thank you

Efficient Mining Algorithms for Large-scale Graphs

NTT Technical Review, Dec 2013, Vol. 11, No. 12

Advanced Processing and Analytics for Large-scale Graphs

NTT Technical Review, Feb 2016, Vol. 14, No. 2

