



電子情報 143

修士論文

リンク解析を用いたウェブ上のスパム発見手法
に関する研究

指導教官 喜連川 優 教授

2006 年2 月3 日提出

情報理工学系研究科 電子情報学専攻

46405 小野 拓史

内容梗概

ウェブスパムとは検索エンジンの検索結果において特定のサイトのランキングを不正に向上させることを目的とした行為のことを指す。ウェブスパムによって検索結果はページの質とは無関係になり精度が低下するため、検索エンジンのインデックスからスパムサイトを除去することが重要な課題となっている。ウェブスパムの手法のなかで、特に関連サイト同士で密にリンクを張ることによってランキングの向上を図るものをリンクスパムと呼ぶ。本論文ではリンクスパムの種類と分布について、日本のウェブアーカイブを用いて調べた結果について考察する。

目次

1	はじめに	1
2	既存の研究	3
3	スパム抽出手法	6
3.1	リンクファーム型スパムの抽出	6
3.1.1	極大クリーク抽出法を利用したリンクスパム抽出	8
3.1.2	近似的極大クリーク抽出法の提案	9
3.2	リンクファーム型スパムとハブ・オーソリティ型スパムを併用したスパムの抽出	11
4	スパムの内容分類	12
5	データセット	15
5.1	ウェブアーカイブ	15
5.2	サイトグラフの作成	18
6	極大クリークの抽出	20
7	近似的極大クリークの抽出	26
8	リンクファーム型スパムとハブ・オーソリティ型スパムを併用したスパムの抽出	33
9	まとめと今後の課題	42
	文献	43
	謝辞	45

目 次

1	大きさ 5 のクリーク	7
2	完全 2 部グラフ	10
3	抽出された集合の完全ハブ	11
4	データセットのインリンク数分布	16
5	データセットのアウトリンク数分布	17
6	サイトグラフ, 無向サイトグラフのインリンク数分布	19
7	クリークのサイズごとに出現するページ数の分布	21
8	クリークのサイズごとに出現するページのドメイン分布	23
9	近似的極大クリークのサイズごとに得られたサイト数の分布	27
10	極大クリーク, 近似的極大クリークのサイズ分布	28
11	近似極大クリークのサイズごとに出現するサイトのドメイン分布	29
12	完全ハブを n 個持つ極大クリークの数	34
13	完全ハブを n 個持つ近似的極大クリークの数	37

表 目 次

1	販売促進サイトの例 1	13
2	販売促進サイトの例 2	14
3	極大クリーク計算のための所要時間	20
4	極大クリークのスパム分類	24
5	jp ドメインの極大クリークの内容	25
6	近似的極大クリークのスパム分類 (小さいサイズ)	30
7	近似的極大クリークのスパム分類 (大きいサイズ)	31
8	jp ドメインの近似的極大クリークの内容	32
9	極大クリークの完全ハブサイトのサンプリング結果	35
10	近似的極大クリークの完全ハブサイトのサンプリング結果	38
11	jp ドメインページを含む極大クリークの割合	39
12	jp ドメインの極大クリークの内容	40
13	完全ハブサイトを持たない近似的極大クリークの割合	41

1 はじめに

今日のウェブにおいて検索は日々の情報収集のために必要不可欠なツールとなっている。このため情報提供者にとって検索エンジンにおける結果の上位に自身のページがランキングされることは重要な意味を持つ。特に通信販売などの商用サイトにおいては、閲覧者を多数獲得する必要があることから、ページの記述方法やサイトの構成方法などを工夫してランキングを改善するSEO (Search Engine Optimization) と呼ばれる手法が多く用いられている。現在のウェブにおいてはSEO 手法を悪用して検索エンジンに意図的に誤った評価をさせ、実際よりも高いランキングを得ることを主目的としてコンテンツを構成する行為が頻繁に行われている。この行為はウェブスパムと呼ばれ、これを行う主体をスパマーと呼ぶ。ウェブスパムが行われると、検索結果に関連のないページが多く含まれるようになり、検索結果に偏りをもたらしたり、得られる情報の品質の劣化を招いたりすることになる。ウェブスパムの手法は、主として二種類の方法に大別できる。一つはページのテキストを検索エンジンのクエリに適合するように調整する手法であり、関連するキーワードを多数ページに付加するなどの手法が用いられる。もう一方はリンク解析を用いたランキングを行う検索エンジンをターゲットとして、自身のサイトの周辺におけるリンク構造を操作してランキングを上げることを目的とするリンクスパム手法である。多数のサイトを作成しその間に密にリンクを張ることによって、参照数を基にしたランキングを欺くなどの手法が用いられている。現在、主要な検索エンジンではリンク解析が重要な要素としてランキングに用いられていることから、リンクスパムは頻繁に行われており、これへの対処が重要な課題となっている。我々は、ウェブ上においてどのようなリンクスパムがどの程度行われているかを調査し、スパム対策手法を開発することを目標としている。その第一歩として、本論文ではリンクスパムを行う際にウェブのグラフ構造上に現れる大きなクリークに着目し、それらの

分布や種類について日本のウェブアーカイブを用いて調査した結果について報告する。本論文の構成は以下のとおりである。第2章では、本研究に関連する既存の研究について述べる。第3章では本研究において用いたスパム抽出手法について説明する。第4章では抽出されたサイトの内容を評価する際のカテゴリ分け方法について説明する。第5章では本研究に使用したウェブアーカイブの詳細と実験に使用したサイトグラフの構成法について述べる。第6章ではウェブアーカイブからの極大クリーク抽出についての実験とその結果について述べる。第7章では近似的極大クリークの抽出実験とその結果について述べる。第8章ではリンクファーム型スパムとハブ・オーソリティ型スパムを併用したスパムの抽出実験の結果を示す。最後に第9章では結論を記す。

2 既存の研究

Henzinger らは [12] においてウェブスパムを検索技術においてもっとも重要な問題の一つにあげている。リンクスパムへの対処に関しては、まず単純な統計を用いたものがあり、Fetterly らはリンクの回数などの統計を用いたウェブスパムの調査を行っている [1]。この研究ではウェブをクロールしたデータセットから得られた統計的な分散を評価した結果、そのうちのいくつかはウェブスパムに密接に関連しているとしている。これは検索エンジンのコーパス内に特定のページの評価を上げるために大量のページを挿入することを目的とした結果生じる現象である。結論として、以下のような性質はスパムページとしての指標になると述べられている。

- ほかにページの URL に含まれるキーワードの内容や、
- 極端に多くのページから参照されている IP アドレス
- インリンク数の分布や、アウトリンク数の分布において極端な値を持っているページ
- コンテンツの過度の複製

リンクスパム手法がターゲットとするリンク解析手法は主に、Page, Brin による PageRank[5]、および、Kleinberg による HITS[7] の2つである [8]。PageRank はあるページに張られているリンクの数は一般的なウェブユーザーにとっての重要度をあらわしているという前提に基づいており、あるページの重要度はそれにリンクを張っているページの重要度によって計算される。このため、特定のページへリンクを集中させるリンクスパム手法が多く用いられている。リンクスパムが PageRank アルゴリズムに及ぼす影響は [2],[3] に述べられている。[2] においてはクリーク型の構造や、クリーク内の接続数が少ないものに対して PageRank の値をシミュレーションにより計算した。ほかにスター型や、リング型などの構造を調べ、どのような構造が効果的なページ

ランクの上昇に結びつくかということを計算している。実際のウェブ上のデータからそのような構造を抽出し、それぞれのページがどのような PageRank を持っているかを計算した。論文では完全部分グラフの中で、張られているリンク数の数に従って計算される PageRank も増大していく傾向が示されている。[3]において、リンクスパムがいかんして PageRank スコアを改変するかが考察されていて、その方法として二通りのものに分けている。まずあるスパマーが特定のページのランキングを上昇させる方法について述べられていて、その次に複数のスパマーのグループが相互に接続されたスパムファームを構成した場合について、考えられている。スパマーは相互の利益のために、あるいは、経済的な契約に基づいて協力する可能性がある。PageRank の改善手法については、Gyongyi らによって TrustRank[9] が提案されている。TrustRank はスパムでないことが判明しているページからスコアを伝播することでスパムサイトへ高いスコアを割り当て難くしている。HITS はウェブ上のすべてのページについてハブスコアとオーソリティスコアを割り当てる。HITS の定義によると、重要なハブページは多くの重要なオーソリティページにリンクしているものであり、また、重要なオーソリティページは多くのハブページからリンクされているものである。HITS アルゴリズムを用いる検索エンジンはもっとも高いハブスコアとオーソリティスコアを持つページをあわせたものを検索結果として返す。ハブスコアは知名度が高いページの多くにリンクをはることによって容易に上げることができる。高いオーソリティスコアを得るのは比較的難しく、重要だと思われるハブから多くのリンクを得なければならない。リンクスパムの手法は Gyongyi らによって [8] にまとめられており、以下のような手法が述べられている。

- 複数のサイトで協力して相互にリンクを交換する、または、自分で複数のドメインを取得しその間に密にリンクを張ることで相互に PageRank およびハブ・オーソリティスコアを上げることができる。リンクファーム

ムと呼ばれる。Wuらは密な相互リンクを抽出することでリンクファームを自動的に検出する手法を提案している [6].

- 一部のウェブディレクトリサービスには、誰もがリンクを登録できるものが存在する。またブログ、Wikiにはコメントなどにリンクを付加して登録できる。これらを利用するとスパマーはターゲットページへ外部からのリンクを加えることができる。ウェブディレクトリは高いPageRankスコアとハブスコアを持っていることが多いので、この手法はターゲットページのPageRankとともにオーソリティスコアもあげることができる。
- 一般的に役に立つ情報を持つページのコピーを用意し、それらがランキングを上げたいターゲットページを指すようにする。コピーしたページが他のページからリンクされると、ターゲットページのランキングを上げることができる。Webでは著名なディレクトリがいくつかあり、例として、Yahoo!ディレクトリ (dir.yahoo.com) や、DMOZ Open Directory (dmoz.org) などがあげられる。これらのディレクトリはトピックごとにコンテンツをまとめており、各トピックに対して、関連するサイトを掲載している。スパマーはこのディレクトリの一部あるいはすべてを複製し、巨大なリンク集を構築するのである。

3 スпам抽出手法

本研究では、リンクファーム、およびウェブディレクトリ等へのリンク登録を利用したリンクスパムを行った際にウェブのグラフ上に現れるクリークに着目してスパム構造の調査を行う。ここで扱うウェブのグラフは、各ウェブサイトをノード、サイト間に張られたリンクをエッジとした有向グラフである。これをサイトグラフと呼ぶ。ページ単位のグラフでは、複数のページにまたがるリンク交換を検出し難くなるため、サイト単位のグラフを用いている。以下に本研究で用いたスパム抽出法を述べる。

3.1 リンクファーム型スパムの抽出

リンクファームを用いたサイト同士はサイトグラフ上で密に結合されることになる。サイトグラフから、2つのサイトの間相互にエッジが張られている場合にのみ方向無しエッジが存在する無向グラフを抽出すると、ほぼすべてのリンクファームはクリークを含むことになる。クリークとはすべてのノードが互いにエッジによって相互連結されている部分グラフを指す。例として図1に大きさ5のクリークを示す。リンクファームを用いたスパムをリンクファーム型スパムと呼ぶ。

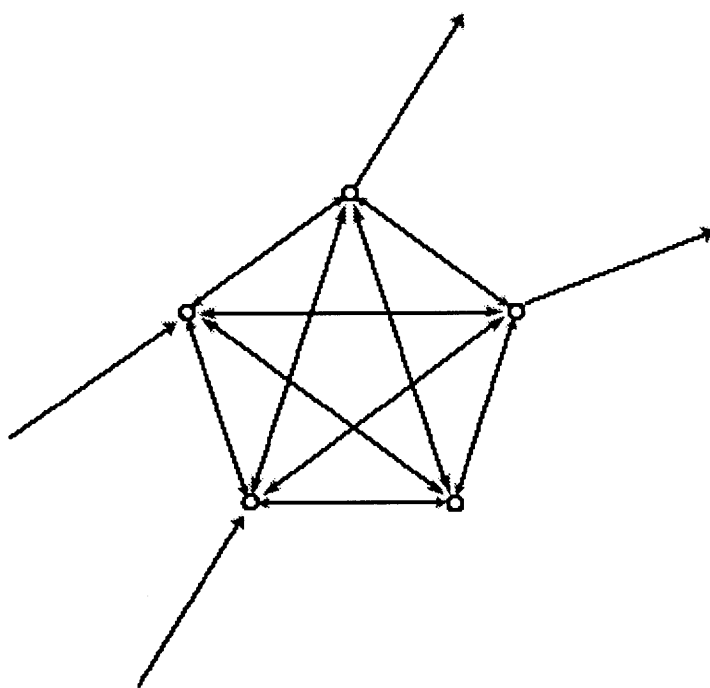


図 1: 大きさ 5 のクリーク

3.1.1 極大クリーク抽出法を利用したリンクスパム抽出

クリークのうち、他のクリークに包含されない極大クリークを抽出することでリンクファームの中心構造を捉えることができる。極大クリーク列挙には、牧野、宇野ら [4] によって提案されたアルゴリズムを使用した。このアルゴリズムはグラフのノード数 n 、エッジ数 m 、最大次数を Δ とすると、極大クリークを一つあたり $O(\Delta^4)$ の計算時間で列挙でき、メモリの使用量は $O(n+m)$ である。このアルゴリズムでは、次数が 80 以上では計算が困難であったので、次数を 80 以下としたサイトグラフにおいて、この実験を行った。第 6 章ではサイトグラフから極大クリークを抽出することでリンクファームを抽出し、その種類について調査する。

3.1.2 近似的極大クリーク抽出法の提案

次に、極大クリークに近い構造を持つ部分を抽出する方法として、共有ノード数によるサイトのクラスタリングを行った。これは以下の手法により行われる。まず、無向グラフの各エッジをソートした状態で読み込む。次に各エッジの両端ノードに対して、双方のノードがリンクしているノードを読み出し、これらのノードの共通要素を調べる。これは最大次数のオーダーで計算可能である。もしこれがある閾値 N 以上ならば同じ集合としてクラスタリングする。クラスタリングのアルゴリズムとしては union-find アルゴリズムを使用する。これは、互いに素な集合を与えられた条件に従ってマージしていくものである。マージされた集合のサイズが N 以上ならばサイトグラフから取り出す。このようにして共通リンク先ノード数が N 以上、サイズ N 以上のクラスタが抽出できる。 N の値を段階的に減少させることで、より共通リンク先ノード数が小さなクラスタが抽出可能である。また、全体の計算量はエッジ数と最大次数の積のオーダーである。第7章では、この手法により得られる極大クリークに類似した形のリンクファームの抽出を行い、その種類について調査する。

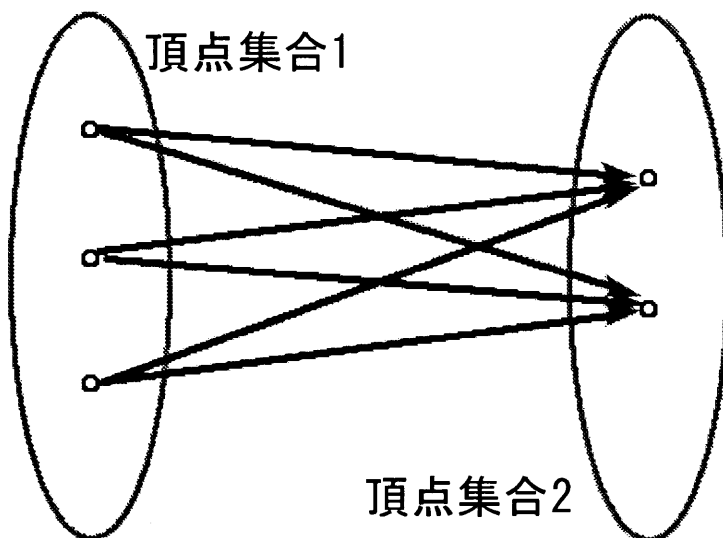


図 2: 完全 2 部グラフ

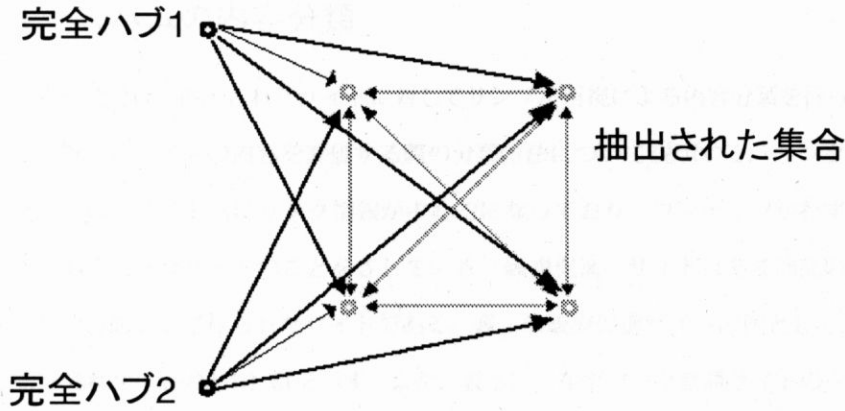


図 3: 抽出された集合の完全ハブ

3.2 リンクファーム型スパムとハブ・オーソリティ型スパムを併用したスパムの抽出

リンク登録を利用したリンクスパムでは、登録サイトの集まりから、ターゲットサイトへの一方向のエッジが密に張られることになる。この構造は、2部クリークを用いて捉えることができる。2つのノード集合を持ち、任意のエッジが2つのノード集合間を結んでいるグラフを2部グラフと呼び、この部分グラフのうち2つのノード集合に含まれるノード同士が全て結合されているものを2部クリークと呼ぶ(図2)。リンク登録を用いたリンクスパムは、登録サイトからターゲットサイトへ完全にエッジが張られた2部クリークを形成することになる。このようなスパムをハブ・オーソリティ型スパムと呼ぶ。ハブ・オーソリティ型スパムの調査手法として、先に得られたリンクファーム型スパムの完全ハブとなるようなサイトを抽出する。完全ハブとは図3のようにある集合のすべてのノードに対してエッジを持っているようなノードを指す。第8章では、極大クリーク、近似的極大クリークの完全ハブを抽出することでハブ・オーソリティ型のスパムを併用しているか調査する。

4 スパムの内容分類

本研究では、抽出されたサイトに対してサンプル目視による内容分類を行った。以下にサイトの内容を評価する際の分類方法について説明する。「リンク集」とは、サイト内にリンク情報が中心的になっており、コンテンツがわずかであるようなサイトのことをさしている。「販売促進」サイトはある商品の宣伝を目的として作られたサイトである。表1に販売促進サイトの例として、旅行者サイトのURL群を示す。また、表2に、在宅ワーク勧誘サイトのタイトル例を示す。これらのサイトは類似した構成のページが大量に作られており、同一の業者によるサイトの集合だと考えられる。そして、互いに密なリンク構造を形成しており、リンクスパムの可能性がある。また、オンラインカジノ、懸賞サービスの勧誘サイトなども販売促進サイトに分類した。本研究では、以上の分類に該当するサイトをスパムサイトと判断した。「一般」サイトとはリンク集、販売促進サイト、アダルトサイトのいずれにも該当しないものである。これには個人のサイト、一般企業のサイト、そして公的機関のサイトなどが含まれる。

表 1: 販売促進サイトの例 1

baltimoremd.areaguides.net
bostonma.areaguides.net
canada.areaguides.net
charlottenc.areaguides.net
columbusoh.areaguides.net
denverco.areaguides.net
detroitmi.areaguides.net
indianapolisin.areaguides.net
international.areaguides.net
jacksonvillefl.areaguides.net
kansascitymo.areaguides.net
lasvegasnv.areaguides.net
longbeachca.areaguides.net
losangelesca.areaguides.net
memphistn.areaguides.net
milwaukeewi.areaguides.net
minneapolismn.areaguides.net
nashvilletn.areaguides.net
neworleansla.areaguides.net
newyorkny.areaguides.net
omahane.areaguides.net
philadelphiapa.areaguides.net
saintlouismo.areaguides.net

表 2: 販売促進サイトの例 2

副業のご案内！安心して取り組めるビジネスです！！
サイドビジネスに最適！！空き時間の有効利用！！
自宅でお仕事しましょう・・！在宅勤務で副収入，副業に最適。
自宅でパソコンを使って出来るビジネスのご紹介です！！
在宅ビジネスの決定版。このビジネスチャンスを Get してください！！
パソコン好きな人へ！あなたの空き時間で副業ができます！！
自宅のPCで空き時間を利用して取り組みます！！
仕事・子育て・家事と両立できる在宅ビジネス！！
副収入への道！サイドビジネスで副収入を得ませんか！！

5 データセット

5.1 ウェブアーカイブ

本実験に使用したデータセットは2004年5月に日本語のウェブページを大規模にクロールをしたものである。クロールの方法としては、jpドメインのページと日本語で書かれた海外ドメインのページを収集する形をとる。jpドメイン以外のサイトについては数ページクロールして日本語のページを発見できない場合、そのサイトのクロールを停止する。アーカイブデータは9600万ページ、45億のリンクから構成されている。クロールの結果得られたデータのインリンク数、アウトリンク数の分布を図4に示した。

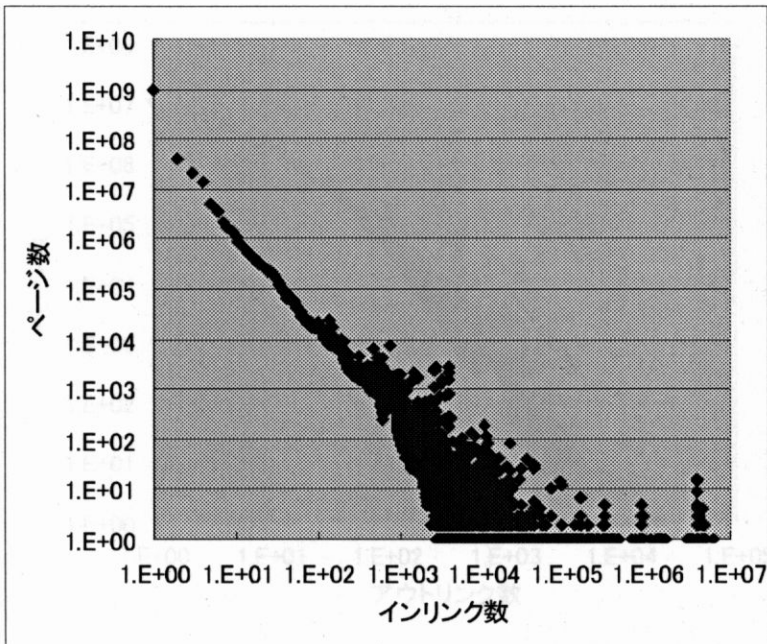


図 4: データセットのインリンク数分布

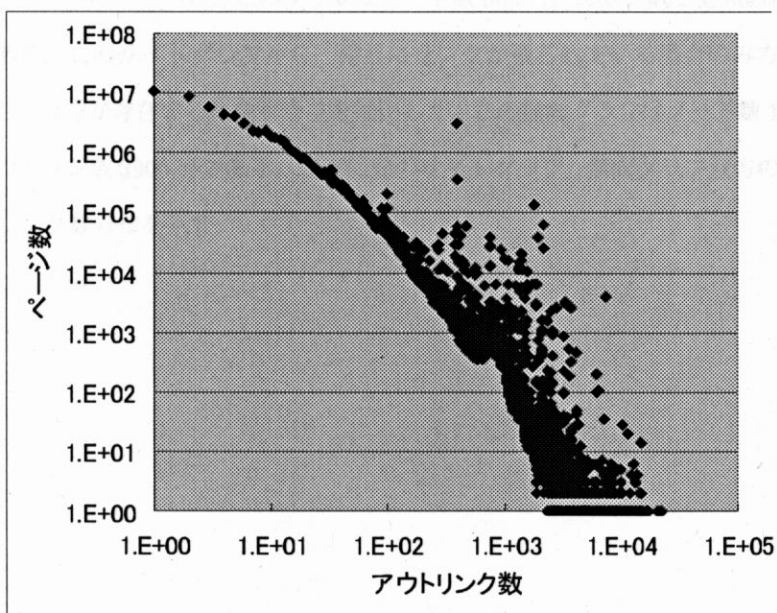


図 5: データセットのアウトリンク数分布

5.2 サイトグラフの作成

データセットからサイトグラフを作成する。このサイトのトップページとしてアーカイブ内にあるインリンク数が3以上のページを集めた。このようなページをシードページと呼び、シードページをサイトのトップページとしたサイトグラフを構成した。各シードページはURLがhttp://A/B/C/の形をしているものに階層を限定した。これは同一の団体によるページを一つにまとめる目的で行った。このグラフはノード数680万、エッジ数2億8000万である。このサイトグラフから、相互にリンクが張られている場合のみ方向無しエッジが存在する無向グラフを抽出した。この無向グラフはノード数160万、エッジ数3900万である。以下に元のサイトグラフ、無向グラフ双方のインリンク数分布を示す。

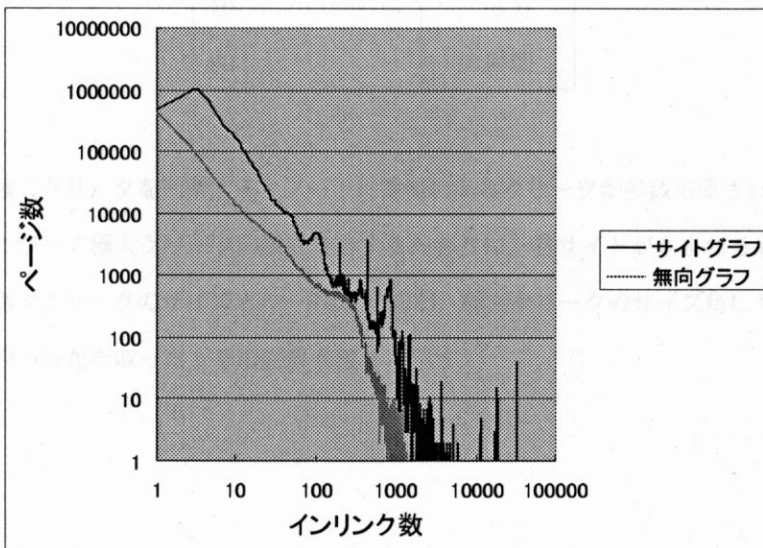


図 6: サイトグラフ, 無向サイトグラフのインリンク数分布

6 極大クリークの抽出

実験は Itanium 2 1.6GHz x 8, メモリー 128GB のマシンで行った. この計算のために所要した時間を表 3 に示した.

表 3: 極大クリーク計算のための所要時間

ノードの最大次数	計算時間
50	1 分
70	6 分
80	18 時間

極大クリークを列挙するとノードに重複のあるクリークが多数生成される. したがって極大クリークの数を集計するかわりに, 各サイトが含まれる最大の極大クリークのサイズ (ノード数) を求め, 極大クリークのサイズ毎にサイト数の分布を取った. その結果を図 7 に示す.

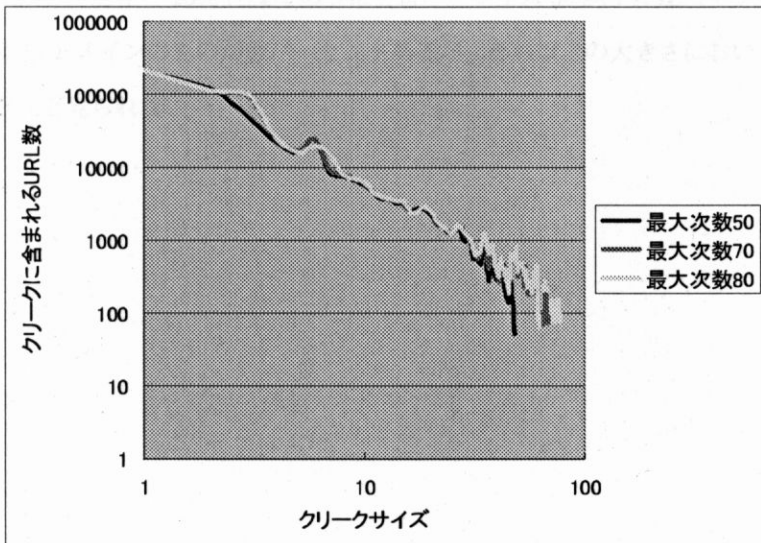


図 7: クリックのサイズごとに出現するページ数の分布

なお、元のサイトグラフとして、各シードページはインリンク数3以上のものとしたため、サイズ3以上のデータを示した。極大クリークのサイズ分布がべき乗則に従うことが分かる。最大次数を80としたとき、サイズ3以上の極大クリークに含まれるノードの総数は60万であり、元のサイトグラフに含まれるノード数の37.5%であった。このときのグラフの極大極大クリークに含まれるサイトの主なドメインの内訳を図8に示す。全体的に国外のサイトが多く、特に.comドメインのものが大半であることがわかる。jpドメインのサイトは全体の16%程度であった。また、サイズ30以下の極大クリークにはjpドメインのものが数パーセントあるが、それ以上の大きさにおいてはほとんど見られない。

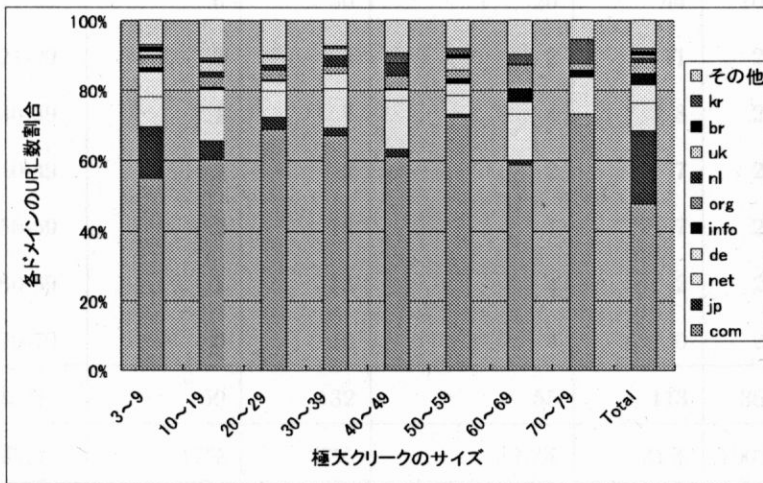


図 8: クリックのサイズごとに出現するページのドメイン分布

次に、極大クリークをランダムに選び、その内容を確認したところ、表4のようであった。全体の83%がスパムサイトと考えられ、アダルトサイトと販売促進サイトの割合が多い。極大クリークのサイズが大きくなるに従い、一般サイトの割合が減少している。

表 4: 極大クリークのスパム分類

サイズ	一般サイト	販売促進	リンク集サイト	アダルト	総計
3-10	38	50	18	34	140
11-20	6	30	20	44	100
21-29	3	4	2	11	20
30-39	1	7	4	8	20
40-49	9	2	2	7	20
50-59	2	10	1	7	20
60-69	1	13	4	2	20
70-79	0	16	4	0	20
総計	60	132	55	113	360
割合	17%	37%	15%	31%	100%

次にjpドメインのページを含む極大クリークを抽出した。これによって得られた極大クリークはサイズが3から29のものであり、その中からランダムに選び、その内容を確認した結果を示したのが表5である。jpドメインでのサンプリング結果ではスパムが5%程度である。アーカイブ全体をサンプリングしたところ、大多数がスパムページであったのに対し、jpドメインでサンプリングを行った場合、スパムページはほとんど見られなかった所が大きく異なった。

表 5: jpドメインの極大クリークの内容

サイズ	一般サイト	販売促進	リンク集サイト	アダルト	総計
3-10	139	0	0	2	141
11-20	51	0	1	2	54
21-29	33	1	2	4	40
総計	223	1	3	8	235
割合	95%	0%	1%	3%	100%

7 近似的極大クリークの抽出

クラスタリングにより得られた集合を近似的極大クリークと呼ぶ。まず、近似的極大クリークのサイズに対して得られたサイト数の分布を図9に示す。これは、べき乗則に従っているといえる。図10では得られた近似的極大クリークのサイズ分布を極大クリークのものとはべて示す。この結果より極大クリークと比較して近似的極大クリークはサイズが大きいことがわかる。その理由として近似的極大クリークは極大クリークと比較して抽出される集合の制約が緩和されているためと考えられる。このときのグラフの近似的極大クリークに含まれるサイトの主なドメインの内訳を図11に示す。極大クリーク抽出のときに得られた結果と同じ傾向にあると言える。

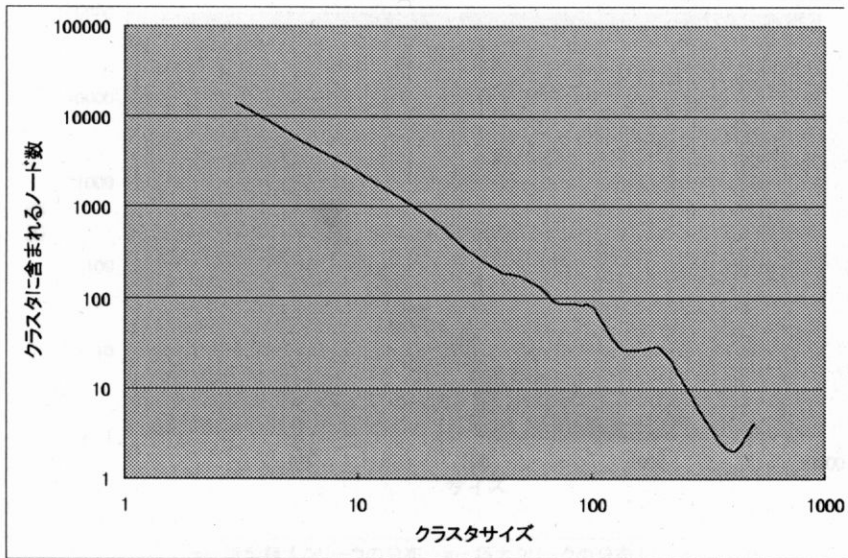


図 9: 近似的極大クリークのサイズごとに得られたサイト数の分布

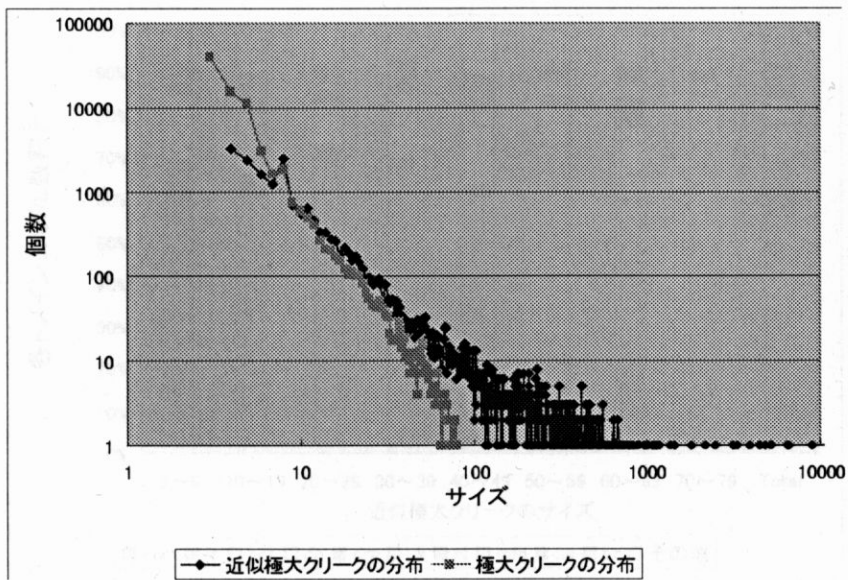


図 10: 極大クリーク, 近似的極大クリークのサイズ分布

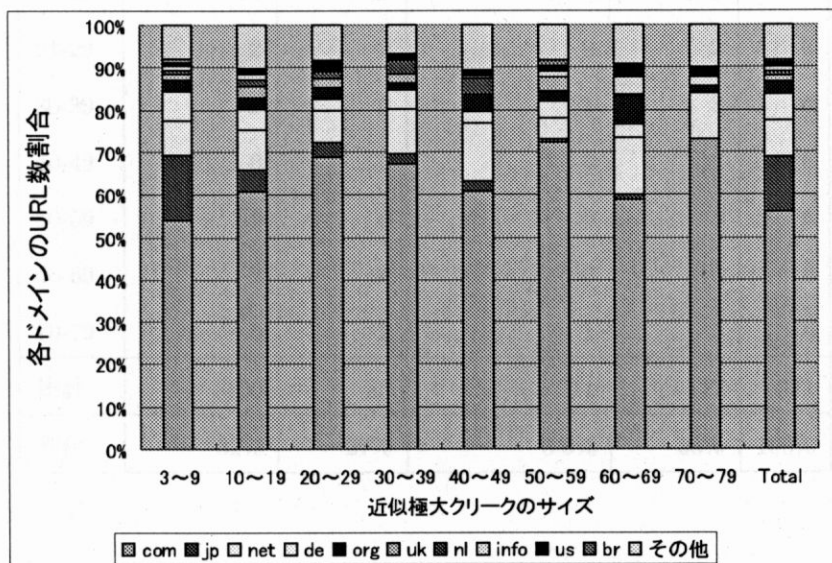


図 11: 近似極大クレークのサイズごとに出現するサイトのドメイン分布

近似的極大クリークを極大クリークと同じ大きさの領域においてランダムに選び、その内容を確認したところ、表6のようであった。スパムの比率は68%程度であった。

表 6: 近似的極大クリークのスパム分類 (小さいサイズ)

サイズ	一般サイト	販売促進	リンク集サイト	アダルト	総計
3-10	78	61	6	92	237
11-20	15	11	3	9	38
21-29	2	2	0	2	6
30-39	2	0	1	3	6
40-49	0	3	0	3	6
50-59	1	0	0	5	6
60-69	2	4	0	0	6
70-79	0	2	1	3	6
総計	100	83	11	117	311
割合	32%	27%	3.5%	38%	100%

また、より大きなサイズの近似的極大クリークについてもランダムに選び、その内容を確認したところ、表 7 が得られた。スパムの比率は 99% 程度であり、そのうちアダルトが大きな割合を占めていることがわかる。

表 7: 近似的極大クリークのスパム分類 (大きいサイズ)

サイズ	一般サイト	販売促進	リンク集サイト	アダルト	総計
100-199	0	5	1	28	34
200-299	1	6	2	17	26
300-	0	0	0	9	9
総計	1	11	3	54	69
割合	1%	16%	4%	78%	100%

次に jp ドメインのページを含む近似的極大クリークを抽出した。これによって得られたものの中からランダムに選び、その内容を確認した結果を示したのが表 8 である。その結果、スパムが占める割合が極大クリークの jp ドメインページにおける割合より大きいことが分かった。

表 8: jp ドメインの近似的極大クリークの内容

サイズ	一般サイト	販売促進	リンク集サイト	アダルト	総計
3 - 10	60	3	0	1	64
11 ~ 20	44	4	0	4	52
21 ~ 29	24	3	0	2	29
30 - 39	17	1	0	2	20
40 - 49	11	1	0	0	12
50 - 59	3	1	0	0	4
60 - 69	7	2	0	1	10
70 - 79	1	2	0	3	6
総計	167	17	0	13	197
割合	85%	9%	0%	7%	100%

8 リンクファーム型スパムとハブ・オーソリティ型 スパムを併用したスパムの抽出

第6章で得られた極大クリークの内、サイズが10以上のものについて完全ハブを持つものの数の分布を示したものが図12である。これより、完全ハブを持たない極大クリークは全体の14%程度であり、半数の極大クリークは4以上の完全ハブを持っていることが分かる。従ってリンクファーム型スパムのほとんどがハブ・オーソリティ型スパムを併用しているといえる。

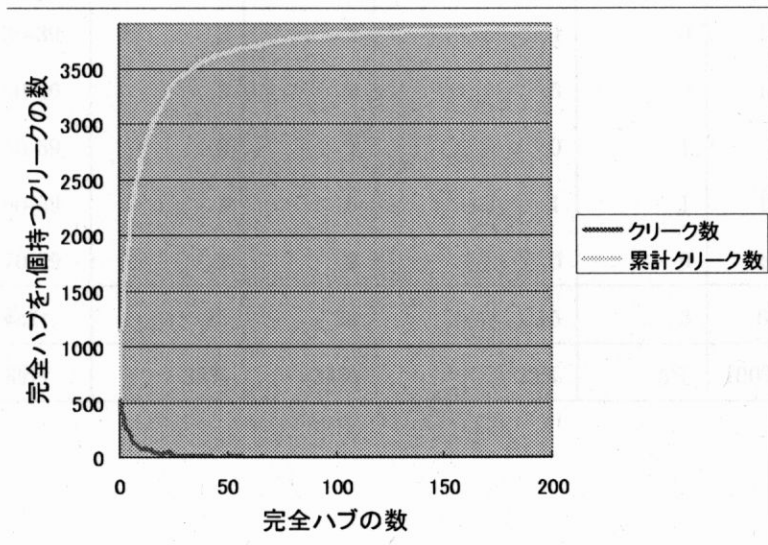


図 12: 完全ハブを n 個持つ極大クリークの数

極大クリークの完全ハブサイトをランダムに選び、その内容を確認したところ、表9が得られた。これより極大クリーク抽出と比較して、スパムサイトの比率が多くなっていることが分かる。

表 9: 極大クリークの完全ハブサイトのサンプリング結果

サイズ	一般サイト	販売促進	リンク集サイト	アダルト	総計
10-19	6	1	2	1	10
20-29	6	1	3	0	10
30-39	3	3	4	0	10
40-49	3	2	5	0	10
50-59	3	1	0	1	5
60-69	3	5	1	1	10
70-79	1	9	0	0	10
総計	25	22	15	3	65
割合	38%	34%	23%	5%	100%

また、第7章で得られた近似的極大クリークの内、サイズが10以上のものについて完全ハブを持つものの数の分布を示したものが図13である。図においては、完全ハブ以外にも、ハブとして条件が緩和されたものについても統計を取った。これより、完全ハブを持たない極大クリークは全体の14%程度であり、半数の極大クリークは4以上の完全ハブを持っていることが分かる。

第6章で得られた極大クリークの内、サイズが10以上のものについて完全ハブを持つものの数の分布を示したものが図12である。従ってリンクファーム型スパムのほとんどがハブ・オーソリティ型スパムを併用しているといえる。

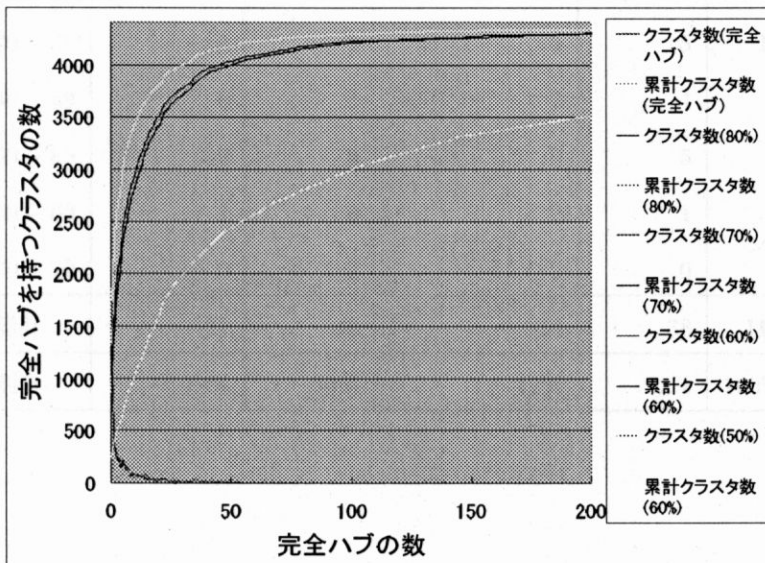


図 13: 完全ハブを n 個持つ近似的極大クリークの数

近似的極大クリークの完全ハブサイトをランダムに選び、その内容を確認したところ、表 10 が得られた。これは近似的極大クリークの内容評価実験と同じ傾向を示している。

表 10: 近似的極大クリークの完全ハブサイトのサンプリング結果

サイズ	一般サイト	販売促進	リンク集サイト	アダルト	総計
11～20	12	22	2	12	48
21～29	7	12	3	11	33
30 - 39	1	7	3	5	16
40 - 49	0	3	2	4	9
50 - 59	0	0	0	5	5
60 - 69	1	0	2	1	4
70 - 79	0	2	2	0	4
総計	21	46	14	38	119
割合	18%	39%	12%	32%	100%

次にこの中から jp ドメインページを含むものを抽出した。表 11 は各サイズにおいて、jp ドメインページを含むものの割合を示している。その一部をランダムに選び、その内容を確認した結果を示したのが表 12 である。スパムの割合が 75% になっており、近似的極大クリークの jp ドメインページにおける割合より増加していることが分かる。

表 11: jp ドメインページを含む極大クリークの割合

サイズ	jp ドメイン	全体	割合
10 - 19	181	3979	5%
20 - 29	16	1242	1%
30 - 39	6	531	1%
40 - 49	0	275	0%
50 - 59	5	228	2%
60 - 69	2	154	1%
70 - 79	0	123	0%
総計	50	6532	1%

表 12: jp ドメインの極大クリークの内容

サイズ	一般サイト	販売促進	リンク集サイト	アダルト	総計
11～20	4	2	1	12	19
21～29	2	5	1	2	10
30 - 39	0	1	1	1	3
40 - 49	0	2	0	0	2
50 - 59	0	0	0	0	0
60 - 69	2	2	0	0	4
70 - 79	2	0	0	0	2
総計	10	12	3	15	40
割合	25%	30%	8%	38%	100%

表 13 では完全ハブサイトを持たない近似的極大クリークの割合を示す。完全ハブを持たない極大クリークは全体の 14% であったが、近似的極大クリークに対しては完全ハブサイトの条件も緩和すれば、極大クリークの結果に近づくことが結果としてわかる。この理由としては、近似的極大クリークの抽出条件も緩和されているので、より大きな集合が抽出されたために完全ハブとなるサイト数が減少するためだと考えられる。

表 13: 完全ハブサイトを持たない近似的極大クリークの割合

ハブサイトの条件	ハブを持たないクリークの割合
完全ハブ	40%
80%	30%
70%	27%
60%	23%
50%	6%

9 まとめと今後の課題

本研究では大規模アーカイブを用いてスパムサイトを抽出しその分布を調べた。そして近似的極大クリークを抽出することにより、巨大スパムを抽出することができた。抽出された極大クリーク型のスパムからサンプルを取り、内容を確認したところ全アーカイブ内から取ったときは83%がスパムであったが、jpドメインについてはスパムページが2割程度と、少ないことが分かった。よってリンクスパムは主に.comなどの国外ドメインにおいて行われており、jpドメインにおいては2004年5月の段階ではリンクスパムがあまり行われていないと言える。また、リンクファーム型とハブ・オーソリティ型の両構造を併用しているサイトは多く、それらはほぼすべてスパムであることも分かった。この手法を用いているサイトを抽出することで、jpドメインサイトにおいて75%の割合でスパムを抽出することが可能であった。

今後の課題としては、異なる抽出条件を用いることによって、本研究では抽出できなかったスパムを発見することや、スパム抽出の精度を向上させることにより、より効率的なスパムの検出を行うことが考えられる。このような手法が確立されればスパムの自動的な判定も可能になり、リンクスパムが検索エンジンに及ぼす悪影響を低減させることができる。

参考文献

- [1] Fetterly, D., Manasse, M., and Najork, M., “Spam, damn spam, and statistics,” Proc. 7th International Workshop on the Web and Databases (WebDB) Paris, France, June 2004.
- [2] Baeza-Yates, R., Castillo, C., and Lopez, V. “PageRank increase under different collusion topologies,” Proc. First International Workshop on Adversarial Information Retrieval on the Web (AIRWeb) Tokyo, Japan, May 2005.
- [3] Gyongyi, Z. and Garcia-Molina, H. “Link spam alliances,” Proc. 31st International Conference on Very Large Data Bases (VLDB) Trondheim, Norway, August 2005.
- [4] Makino, K., and Uno, T., “New Algorithms for Enumerating All Maximal Cliques,” T. SWAT 2004, LNCS 3111, pp. 260-272, 2004.
- [5] Page, L., Brin, S., Motwani, R., and Winograd, T., “The PageRank citation ranking: Bringing order to the web,” Tech. rep., Stanford University, 1998.
- [6] Wu, B., and Davison, B., “Identifying link farm pages,” Proc. 14th International World Wide Web Conference (WWW), Tokyo, Japan, May 2005.
- [7] Kleinberg, J., “Authoritative sources in a hyperlinked environment,” Journal of the ACM, 46(5),
- [8] Gyongyi, Z. and Garcia-Molina, H. “Web Spam Taxonomy,” Proc. First International Workshop on Adversarial Information Retrieval on the Web (AIRWeb) Tokyo, Japan, May 2005.

- [9] Gyongyi, Z., Garcia-Molina, H., and Pedersen, J. “Combating web spam with TrustRank,” In Proceedings of the 30th International Conference on Very Large Data Bases (VLDB) Toronto, Canada, August 2004.
- [10] Benczur, A., Csalogany, K., Sarlos, T., and Uher, M., “SpamRank - Fully Automatic Link Spam Detection,” Proc. First International Workshop on Adversarial Information Retrieval on the Web (AIRWeb) Tokyo, Japan, May 2005.
- [11] Bifet, A., Castillo, C., Chirita, P., and Weber, I., “An Analysis of Factors Used in Search Engine Ranking,” Proc. First International Workshop on Adversarial Information Retrieval on the Web (AIRWeb) Tokyo, Japan, May 2005.
- [12] Henzinger M., Motwani, R., and Silverstein, C. Challenges in Web Search Engines. SIGIR Forum 36(2), 2002.

謝辞

本研究を進めるにあたって、指導教官の喜連川優教授より非常に興味深い研究テーマを頂くと共に研究全般に渡って懇切丁寧な指導を頂きました。非常に感謝しております。

また、豊田正史先生にも日頃から丁寧に指導を頂き、とても感謝しております。さらに、日頃お世話になった研究室の皆様にもこの場を借りて深く感謝の意を表します。