

最適なウェブサイト構造の調査

株式会社Faber Company (配属先: Faber & Technology)

B143369 廣中 詩織

ウェブサイト内の ページ数が増える

- ・重複ページを作成してしまう
 - ・サイト構造の把握が難しい

- ① 「サイト内ページリストの作成」
 - ② 「リンクネットワークの作成と分析」
 - ③ 「リンクネットワークの可視化」

① サイト内ページリストの作成

目的 重複ページやタイトルの設定ミスなどを発見しやすくする
→クローラがサイト内を巡回してページをリスト化

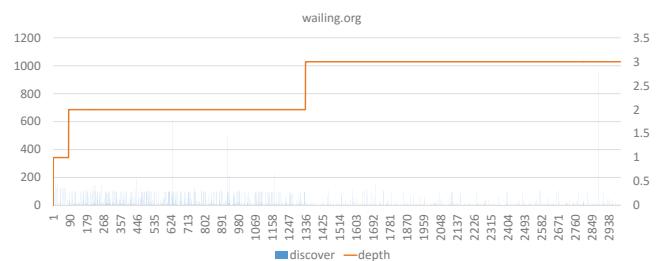
問題 ページ（コンテンツ）の取得をいつまで続けるか？

CGIなどによってURLを自動生成するサイトが存在する

→ トップページからリンクをたどる深さを設定 + 時間制限

表: サイト内のページ一覧(抜粋)。ページごとに設定されているタイトルや見出し、キーワードなどのパラメータを一覧で確認できる。

URL	depth	status_code	content_type	size	title	h1	h2	keywords	description	canonical
.../archives/3853?replaytocom=73	3	200	text/html; charset=UTF-8	40148	約7万円で購入できる3Dプリンタ『3D...』	3DプリンターのJapanese Makers, 約7万円で購入できる3D... 7万円で購入できる3D...	3Dプリンタ『3D...』 安,格安,3dプリンタ	3Dプリンタ,メンテル型,激メンテル型の激安3D... 3Dプリンタが秋...	.../archives/3853	
.../archives/tag/reprap	3	200	text/html; charset=UTF-8	22724	Reprap 3DプリンターのJapanese Makers	3DプリンターのJapanese Makers,Reprap	約7万円で購入できる3D... 3Dプリンタ『3D...』 安,格安	約7万円で購入できる3Dプリンタ,メンテル型,激... 3Dプリンタ『3D...』 安,格安	.../archives/tag/reprap	
.../archives/958	1	200	text/html; charset=UTF-8	39994	Arduinoではんだ付け無しで簡単な回路を作る3DプリンターのJapanese Makersは...んだけ付け無しで簡単... 3Dプリンターの比...	Arduinoではんだ付け無しで簡単な回路を作る3D... んだけ付け無しで簡単... 3Dデーター共有サイトより3Dデータをダウン	フレッドボード ジャンバ ブレッドボードと... 單な回路を作る ... ~ ワイヤー,arduino,回路 3Dデーター共有サイト blade 1,thingiverse,3Dブリ	フレッドボード ジャンバ ブレッドボードと... 單な回路を作る ... ~ ワイヤー,arduino,回路 3Dデーター共有サイト blade 1,thingiverse,3Dブリ	フレッドボード ジャンバ ブレッドボードと... 單な回路を作る ... ~ ワイヤー,arduino,回路 3Dデーター共有サイト blade 1,thingiverse,3Dブリ	.../archives/958
.../archives/3106?replaytocom=35	3	200	text/html; charset=UTF-8	48335	3Dデーター共有サイトより3Dデータをダウンロードして...	3DプリンターのJapanese ...	3Dデーター共有サイト ... より3Dデータを ...	3Dデーター共有サイト ... より3Dデータを ...	3Dデーター共有サイト ... より3Dデータを/archives/3106
.../archives/1689?replaytocom=90	3	200	text/html; charset=UTF-8	44167	家庭用3Dプリンターを見る場所 3D... 家庭用3Dプリンタ... 見学,え見える,体験,展示,デ... 3D... 検討している...	3DプリンターのJapanese Makers,家庭用3Dプリンタ... 見学,え見える,体験,展示,デ... 3D... 検討している...	3Dプリンターを見る場所 3D... 家庭用3Dプリンタ... 見学,え見える,体験,展示,デ... 3D... 検討している...	3Dプリンターを見る場所 3D... 家庭用3Dプリンタ... 見学,え見える,体験,展示,デ... 3D... 検討している...	.../archives/1689	



図：クローリング中の新規発見リンク数の推移

② リンクネットワークの作成と分析

目的 問題のあるコンテンツを発見する

見せたいページ（コンテンツ）がちゃんとリンクされているか？

- 取得したページ間のリンクからリンクネットワークを作成
 - リンクネットワークのリンク構造から計算したPageRankとアクセス数から計算したページ訪問率を比較

PageRank

= ランダムにリンクをたどったときそのページにたどり着く確率

表: PageRankとサイト内部からの
ページアクセス数とを比較したスコアから
リンクの良さをはかることができる可能性がある

page	score
/3d-printer-list	-0.06161
/archives/4105	-0.05593
/archives/2201	-0.05386
/archives/270	-0.03985
/archives/633	-0.03747
/history/サービス	0.01909
/history/動画	0.018975
/history/word	0.018894
/history/管理人のつぶやき	0.018054
/contact	0.017185

③ リンクネットワークの可視化

目的 サイトの構造を視覚的に把握する

ページの共通部分
(メニューなど) から
リンクされているページ

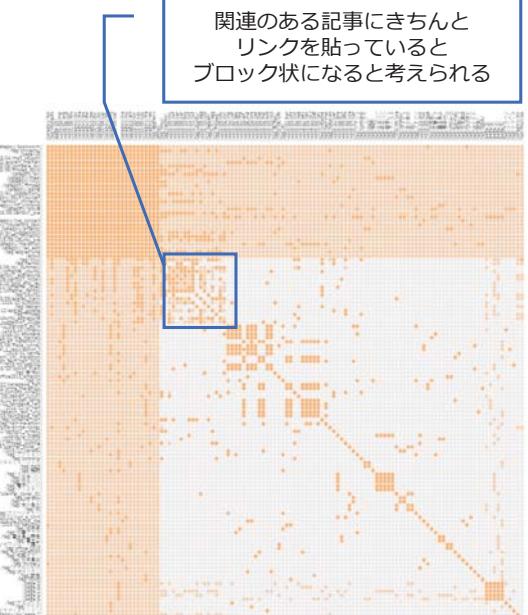


(2) フィルムが乱雑なサイト

図: コンテンツが乱雑なサイトの
ネットワーク図（ノードの大きさ
はPageRankの値をあらわす）

（はPageRankの値をあらわす）
リンクネットワークは密なグラフ
であるため、ネットワーク図によ
る可視化は適していない。

関連のある記事にきちんと
リンクを貼っていると
ブロック状になると考えられる



(b): 構造に少し気を使っているサイト

図: ページ数が同程度のブログサイトから作成したリンクネットワークの隣接行列 (PageRankが高い順に並びかえ) 上側と左側の帯は、ページの共通部分にあるリンクである。また、その他の部分をみると、(a) はランダムにリンクが散らばっているが、(b) はクラスタのようなものができている。