

リンク不整合検出による Web サイト診断 — 論理的な不整合の自動判定

河合英紀 河野泉 石黒義英 福島俊一

NEC インターネットシステム研究所

1. はじめに

近年、企業 Web サイトの大規模化・複雑化に伴い、リンク不整合の問題が深刻化している。リンク不整合には、デッドリンクのように文書にアクセスした時点でエラーが発生する不整合(物理的不整合)と、明示的なエラーは発生しないが論理的な誤りを生じている不整合(論理的な不整合)の2種類がある。

筆者らは、論理的な不整合の検出方法として、複数ページ間でのリンク元、リンク先、アンカー文字列の同一性と仲間外れに着目した検出ルールを提案し、実際の企業 Web サイトでは、物理的不整合に加え論理的な不整合も多数存在することを示してきた[1]。そのなかで、単純な検出ルールでは、論理的な不整合の候補として検出されたリンクが必ずしも不整合でない場合もあるため、候補を人手でチェックする作業の効率向上が望まれていた。そこで本稿では、論理的な不整合の候補に対して、機械学習の一手法である Support Vector Machine(SVM) [2]を用いてそのリンクの整合性を自動判定する方法を提案する。

2. 論理的な不整合検出方法とその問題点

図1に、論理的な不整合の例を示す。図1では、Page C、D、Eからアンカー文字列「Y4100」で、製品 Y4100の詳細情報である Page Aへ正しくリンクが張られている。ところが、Page Fからはアンカー文字列「Y4100」で製品 X3100の詳細情報である Page Bへ間違っただリンクが張られている。ここで、検出ルール(1)「アンカー文字列が同一のリンクのグループ内で、リンク先文書が異なるサブグループ」を適用すると、Page C～Fに含まれるアンカー文字列が「Y4100」の4つのリンクのグループ内で、Page Aをリンク先とする3つのリンクと、Page Bをリンク先とする1つのリンクが異なるサブグループとなり、論理的な不整合の候補として検出される。これら2つのサブグループに対して、それぞれアンカー文字列とリンク先の整合性を人手でチェックすることにより、Page Fからのリンクが論理的な不整合であると判断される。文献[1]では検出ルール(1)の他に4つの検出ルールを用いた。

一方、上記の検出ルールは非常に単純であるため、論理的な不整合の候補として検出されたリンクが必ずしも論理的な不整合でない場合もある。例えば、アン

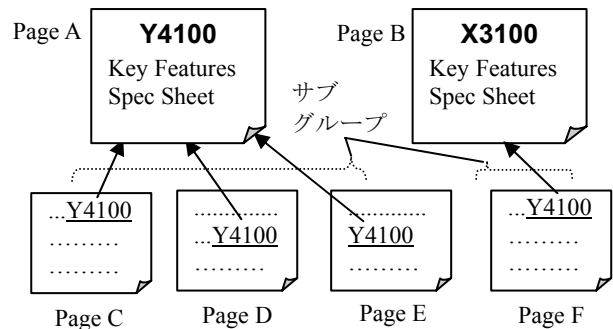


図1 論理的な不整合の例

カー文字列が「戻る」「次へ」などのナビゲーション用リンクは、同じアンカー文字列で異なる文書にリンクされる傾向があるため、検出ルール(1)で誤検出されてしまう。誤検出を削減するためには、検出ルールをさらに詳細化する方法が考えられる。しかし、実際の論理的な不整合のケースは多様であり、全てのルールを網羅的に記述するのは困難である。そこで本稿では、機械学習の一手法である SVMを用いてリンクの整合性を自動判定することにした。

3. SVMによる論理的な不整合の自動判定

SVM[2]は「マージン最大化」を基準として分離平面を決定する2値分類器であり、他の機械学習の手法に比べて優れた認識性能を発揮できるとして、近年注目を集めている。

リンクの整合性判定は、診断対象とするリンクの素性を入力とし、そのリンクが不整合か否かを出力とする2値分類問題と考えることができる。本稿では、検出ルールに該当したリンクのサブグループを診断対象とし、リンクの素性として次の2種類を用いることとした。

(1) リンクグループに関する素性

図1のような論理的な不整合が存在する場合、検出ルール(1)でグループ化するとサブグループに含まれるリンク数に偏りが出る。この時、少数派のサブグループに属するリンクが不整合の可能性が高いと考えられる。そこで、(1a)グループ化に使った検出ルール、(1b)グループ内における各サブグループのリンク数の分散、(1c)リンクグループのリンク数に対するサブグループのリンク数の比を素性とした。

(2) リンク関係に関する素性

間違いリンクなどの論理的な不整合では、アンカー文字列とリンク先文書で使われている単語が一貫し

ない傾向があると考えられる。そこで、各リンクサブグループについて、アンカー文字列中の単語がリンク先文書において、(2a)タイトルに含まれる割合、(2b)本文に含まれる割合、(2c)HタグやFONTタグ等によって強調表示されている割合を素性とした。

リンクの整合性を判定するための素性としては、他にもリンク先 URL に含まれるディレクトリ名やアンカー文字列に含まれる単語なども考えられる。しかし、予備実験では上記2種類の素性に限定した方が、アンカー文字列やURLを素性として追加するよりも優れた判別精度を示した。

また、学習データとしては、(A)診断対象データと同じサイトの過去のチェック済みデータ、(B)診断対象データとは異なるサイトの過去のチェック済みデータ、(C)診断対象データの一部を手でチェックしたデータ、などの利用が考えられる。本稿では、同じサイトを定期的にチェックする運用形態を想定し、(A)診断対象データと同じサイトの過去のチェック済みデータを学習データとして用いることにした。

4. 実験

実際に2件の企業Webサイトを対象として不整合チェック作業の効率改善実験を行った。サイトAに関しては2003年8月15日に不整合チェックを行った結果を学習データとし、2003年9月15日にWebブラウザで収集したデータを診断対象データとした。サイトBも同様に2003年8月28日のチェック済みデータを学習データ、2003年9月26日のデータを診断対象データとした。SVMプログラムには、TinySVM¹を用いた。

5. 結果および考察

表1に検出ルール単独とSVMを用いた場合の判定精度を示す。検出ルール単独で見つかった不整合総数に対し、SVMが検出できた不整合の割合(再現率)は50%台であった。しかし、不整合の候補として検出され実際に不整合だった割合(適合率)は検出ルール単独では20%台だがSVMでは80%台に向上した。

図2および図3にそれぞれ、サイトA、Bにおいて作業時間あたりに検出された不整合数を示す。図2の7月15日、8月15日、および図3の8月28日はSVMを用いずに検出ルールに該当したリンクサブグループを先頭から順にチェックした。一方、図2の9月15日と図3の9月26日はSVMを用いて、分離平面からの距離が大きいサブグループから順に検出ルールに該当したリンクを全てチェックした。いずれの場合においても、SVMを使わない場合は作業時間に従って不整合数はほぼ線形に増加したのに対し、SVMを使った場合は、作業の初期段階で多くの不整合が見つかり、やがて飽和する傾向が見られた。こ

表1 提案手法の精度評価

サイト	検出ルール	SVM	
	適合率	適合率	再現率
A	22% (153/684)	87% (83/95)	54% (83/153)
B	26% (229/881)	80% (125/156)	55% (125/229)

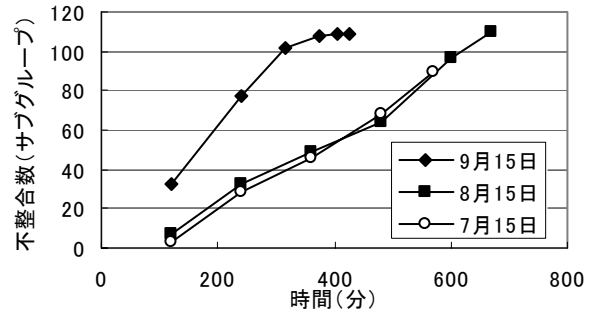


図2 サイトAにおける作業効率

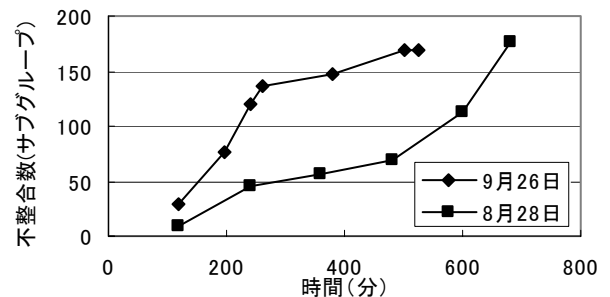


図3 サイトBにおける作業効率

れは、SVMの適合率が高いために、検出ルール単独の場合に比べて短時間で多くの不整合を検出できていることを示唆している。不整合総数の90%を発見するまでにかかる時間を比べると、サイトAでは8月15日に約600分かかったのに対して、9月15日には約300分とほぼ半分に短縮でき、サイトBでも約6割に短縮できた。

6. おわりに

SVMを用いた論理的な不整合の自動判定法を提案し、実際の企業Webサイト上でその有効性を検証した。その結果、再現率は50%台に低下したが、適合率は80%台に向上し、不整合発見の時間をほぼ半分に短縮することができた。提案手法は、Webサイトをくまなく診断する用途には向かないが、短時間で効率よくサイトの問題点を発見する用途に適している。

今後は、自動判定の再現率の向上を目指すとともに、他の運用形態も考慮した学習データの選び方について検討していきたい。

参考文献

- [1] 河合英紀, 河野泉, 石黒義英, 福島俊一, リンク不整合検出によるWebサイト品質評価, FIT2003, LD-005, 2003.
- [2] Vladimir N. Vapnik, Statistical Learning Theory, A Wiley-Interscience Publication, 1998.

¹ <http://cl.aist-nara.ac.jp/~taku-ku/software/TinySVM/>