

ネットワークコンピューティング(1) 情報検索

関西学院大学工学部情報科学科
北村泰彦

1

情報検索とは

- 大量の情報の中から、ユーザの要求を満たす情報を見つけ出すこと。
- ここでは、検索対象となる情報は、 (テキスト)を想定する。
- 検索の対象となる複数の文書を (document collection), ユーザの情報要求を (query)と呼ぶ。

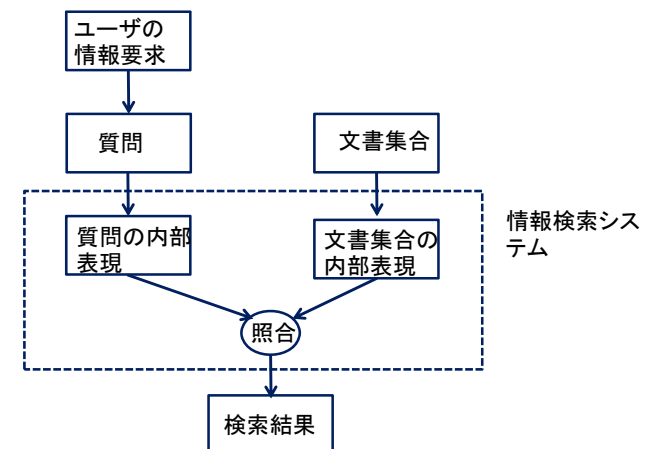
2

情報検索とは

- 代表的な情報検索システムにがある。
 - 文書集合: インターネット上に存在するWebページ
 - 質問: ユーザの与えるキーワード
 - 機能: キーワードを含むWebページを探し出す。

3

情報検索のモデル



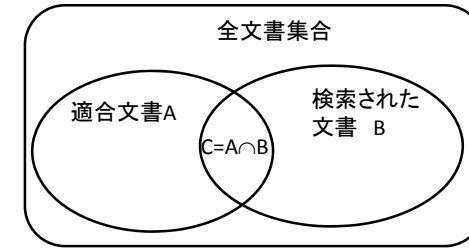
4

データベースとの違い

- データベースシステムは情報検索システムである。
 - 文書集合: データベース
 - 質問: SQL文などの検索式
 - 機能: 検索式を満たすデータをデータベースより抽出する
- ただし, データベースで扱うデータは一般的に (複数の属性値から構成されている) されている. 情報検索システムで扱うデータは文書 (テキスト) であり,

5

情報検索の評価尺度



- (recall): 検索の完全性を評価するための尺度で, 適合文書のうち, 検索された文書の割合を示す. 検索漏れの少なさを示す尺度.
- あるいは (precision): 検索の正確性を評価するための尺度で, 検索された文書のうち, 適合文書の割合を示す. 検索ノイズの少なさを示す尺度.

6

情報検索の評価尺度

- 10000件の文書集合の中に, 関西学院大学に関する文書は100件あるとする. 検索質問文「関西学院大学」で検索したところ, 200件の文書が得られたが, その中で関西学院大学に関する文書は50件であった.
- その時の再現率は である.
- その時の適合率は である.

7

Web情報検索の歴史

- 第一世代: 人手による情報収集
- 第二世代: による情報収集
- 第三世代: による選別強化

8

第一世代

- 人海戦術: ウェブページを選別収集し、それを内容に応じてカテゴリに分類し、 (データベース)に登録する。
- 「Yahoo!カテゴリ」(<http://dir.yahoo.co.jp/>)やOpenDirectory(<http://www.dmoz.org>)が代表例。他の多くの検索エンジンはサービスを終了している。
- 人間の判断に基づいて選りすぐりのページを集めることができる反面、労力の限界から、規模拡大や情報鮮度の維持が難しい。

9

ディレクトリ

- 人手で作った、ウェブページの索引集



第二世代

- (Crawler): インターネットを巡回して、ウェブページを自動収集するシステム。, とも呼ばれる。
- 当初のLycos, AltaVista, gooなどが代表例
- 検索エンジンの規模、カバー率は拡大されたが、検索キーワードに多数のページがヒットし、必要な情報が埋もれてしまう。情報検索のは向上したが、は低下した。

11

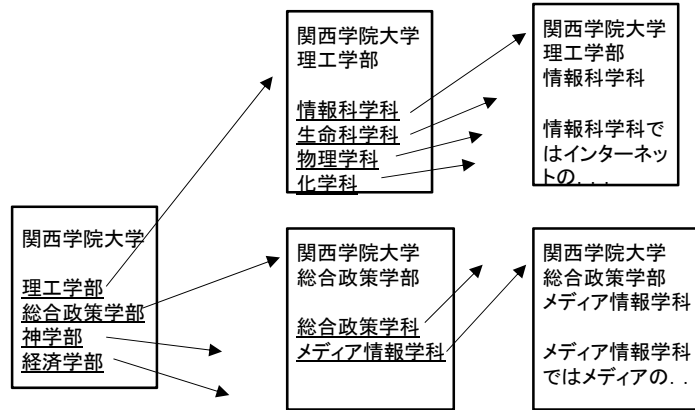
Web情報検索の手順

- Webページの収集: クローラ
- 索引語の抽出
 - 形態素解析
 - 不要語の除去
- 索引語の重み付け
 - 局所的重み付け: 索引語頻度(term frequency: TF)
 - 大域的重み付け: 文書頻度(document frequency: DF)
- 文書検索
 - ベクトル空間モデル

12




ウェブページ中のそれぞれのリンクに対して、その先のページ収集を再帰的に行う。




- 文書を構成する文字列を単語に分割し、各単語に品詞や語形変化などの情報を与える処理。

検索	ケンサク	名詞-サ変接続
エンジン	エンジン	名詞-一般
の	ノ	助詞-連体化
仕組み	シクミ	名詞-一般
と	ト	助詞-並立助詞
技術	ギジュツ	名詞-一般
の	ノ	助詞-連体化
発展	ハツテン	名詞-サ変接続

不要語の除去

-  (stop word): 文書を特定する能力が低く、索引語として適当でない単語。
 - 日本語の助詞: 「は」、「が」など
 - 英語の冠詞: “a”, “the” など
 - 英語の前置詞: “at”, “of”, “in” など

索引語の重み付け

- 索引語の中には文書の内容と密接に関係したものもあるし、そうでないものも存在する。文書の内容を示す上での重要度を表す重み付けを検索語に行うと精度の高い検索が可能になる。
- 局所的重み: l_{ij}
 - 索引語 w_i の文書 d_j における出現頻度に基づき計算される重み。文書中に頻繁に出現する索引語に大きな値が与えられる。
 - 例:  TF: f_{ij} (索引語 w_i の文書 d_j における出現頻度)

索引語の重み付け

- 大域的重み: g_i
 - 文書集合全体における索引語 w_i の分布を考慮して決定される重み. 特定の文書に集中して出現する索引語に対して大きな値が与えられる.
 - 例: の逆数IDF: $\log \frac{n}{n_i}$ ただし n は文書の総数, n_i は索引語 w_i を含む文書数. 対数化するのはIDFの値の変化を小さくするため.

17



- 検索対象となる を d_1, d_2, \dots, d_n とする. これらの文書集合に対する を w_1, w_2, \dots, w_m とする. このとき文書 d_j を以下の で表現される.

$$d_j = \begin{bmatrix} d_{1j} \\ d_{2j} \\ \vdots \\ d_{mj} \end{bmatrix}$$

ここで d_{ij} は索引語 w_i の文書 d_j における重みである.

18

ベクトル空間モデル

- 文書集合全体は $m \times n$ の D によって表現できる.

$$D = [d_1 \quad d_2 \quad \cdots \quad d_n]$$

$$= \begin{bmatrix} d_{11} & d_{12} & \cdots & d_{1n} \\ d_{21} & d_{22} & \cdots & d_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ d_{m1} & d_{m2} & \cdots & d_{mn} \end{bmatrix}$$

索引語・文書行列の各列は であり, 各行は である.

19

ベクトル空間モデル

- 検索質問文に含まれる索引語 w_i の重みを q_i とすると, q は以下のように表される.

$$q = \begin{bmatrix} q_1 \\ q_2 \\ \vdots \\ q_m \end{bmatrix}$$

20

ベクトル空間モデル

- 文書検索は検索質問ベクトル q と各文書ベクトル d の間の $\cos(\mathbf{d}_j, \mathbf{q})$ を計算することで行う。類似度の計算にはコサイン尺度や内積が用いられる。

- コサイン尺度
$$\cos(\mathbf{d}_j, \mathbf{q}) = \frac{\mathbf{d}_j \cdot \mathbf{q}}{\|\mathbf{d}_j\| \|\mathbf{q}\|} = \frac{\sum_{i=1}^m d_{ij} q_i}{\sqrt{\sum_{i=1}^m d_{ij}^2} \sqrt{\sum_{i=1}^m q_i^2}}$$

- 内積

$$\mathbf{d}_j \cdot \mathbf{q} = \sum_{i=1}^m d_{ij} q_i$$

21

ベクトル空間モデル

- 索引語

- w_1 : Bioinformatics
- w_2 : Biology
- w_3 : Chemistry
- w_4 : Enzymes
- w_5 : Evolution
- w_6 : Genes
- w_7 : Genome(s)
- w_8 : Proteins

22

ベクトル空間モデル

- 文書

- d_1 : Bioinformatics: A Practical Guide to the Analysis of Genes and Proteins
- d_2 : Proteins, Enzymes, Genes: The Interplay of Chemistry and Biology
- d_3 : Adaptive Evolution of Genes and Genomes
- d_4 : Advances in Genome Biology: Genes and Genomes
- d_5 : Bioinformatics and Genome Research
- d_6 : Data Analysis in Molecular Biology and Evolution

23

ベクトル空間モデル

- 索引語・文書行列 (索引語の重みは頻度)

$$D = \begin{bmatrix} 1 & 0 & 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 & 0 & 1 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 1 \\ 1 & 1 & 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 2 & 1 & 0 \\ 1 & 1 & 0 & 0 & 0 & 0 \end{bmatrix}$$

24

ベクトル空間モデル

- 検索質問文: Genes and Genomes

$$q = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ 1 \\ 0 \end{bmatrix}$$

25

ベクトル空間モデル

- コサイン尺度を用いた場合の類似度計算

$$\cos(d_1, q) = \frac{1}{\sqrt{3}\sqrt{2}} = 0.408$$

$$\cos(d_2, q) = \frac{1}{\sqrt{5}\sqrt{2}} = 0.316$$

$$\cos(d_3, q) = \frac{2}{\sqrt{3}\sqrt{2}} = 0.816$$

$$\cos(d_4, q) = \frac{3}{\sqrt{6}\sqrt{2}} = 0.866$$

$$\cos(d_5, q) = \frac{1}{\sqrt{2}\sqrt{2}} = 0.5$$

$$\cos(d_6, q) = \frac{0}{\sqrt{2}\sqrt{2}} = 0$$

26

第三世代

- 第二世代の検索エンジンは検索精度が十分でなくなってきた。
 - ウェブページの数が増大になった
 - ほとんどの検索語は1, 2語であり、適合度を計算するための情報に乏しい。
 - ランキングの上位に位置することを意図した行為が存在する
- に基づく検索結果ランキング方式の導入
- Googleが代表例

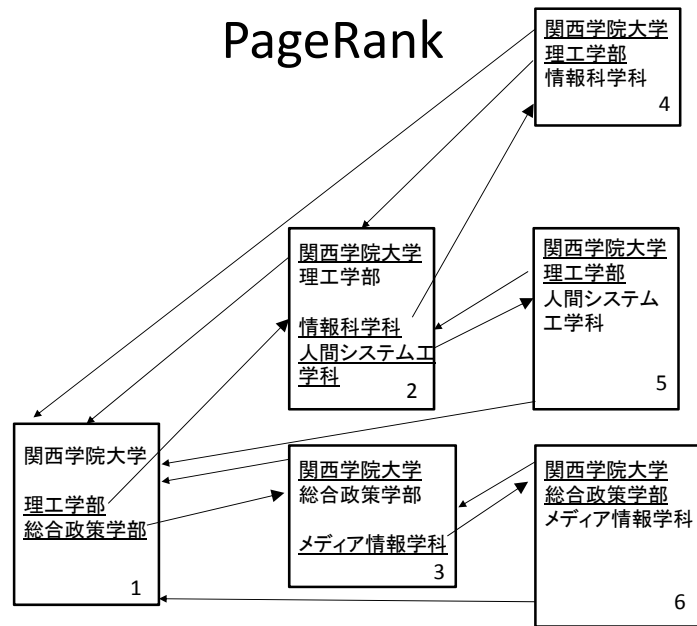
27



- 「多数引用されているページは信頼できる。また、信頼できるページに引用されるページも信頼できる。」という考えに基づくスコア計算方式。
- PageRankは下式を繰り返すことにより求められる。
$$R(p) = \frac{\epsilon}{n} + (1 - \epsilon) \cdot \sum_{(q,p) \in G} \frac{R(q)}{\text{outdegree}(q)}$$
- $R(p)$ はページ p のPageRank, n は対象とする(Webページをノードとし、それらのリンクをエッジとした)グラフ G のノード総数(Webページ数), $\text{outdegree}(q)$ はページ q からの外向きリンク数。また $\sum_{p \in G} R(p) = 1$ であり, ϵ はdampening factorとよばれ, 0.1~0.2の値を取る。
- ユーザは $(1-\epsilon)$ の確率で現在のWebページからのリンクをたどり, ϵ の確率でまったく無関係なWebページにジャンプする。

28

PageRank



PageRank

- 初期値: $R(1)=R(2)=R(3)=R(4)=R(5)=R(6)=1/6$
- 1回目 $\epsilon=0.1$ とする
 - $R(1) = \frac{0.1}{6} + 0.9 \cdot \left(\frac{1}{3} + \frac{1}{2} + \frac{1}{2} + \frac{1}{2} + \frac{1}{2} \right) = 0.367$
 - $R(2) = \frac{0.1}{6} + 0.9 \cdot \left(\frac{1}{2} + \frac{1}{2} + \frac{1}{2} \right) = 0.242$
 - $R(3) = \frac{0.1}{6} + 0.9 \cdot \left(\frac{1}{2} + \frac{1}{2} \right) = 0.167$
 - $R(4) = R(5) = \frac{0.1}{6} + 0.9 \cdot \left(\frac{1}{3} \right) = 0.067$
 - $R(6) = \frac{0.1}{6} + 0.9 \cdot \left(\frac{1}{2} \right) = 0.092$
- 得られた値を初期値として、この計算を繰り返す。

30

参考文献

- 北研二, 津田和彦, 獅々堀正幹: 情報検索アルゴリズム, 共立出版, 2002
- 福島俊一: 検索エンジンの仕組みと技術の発展, 情報の科学と技術, 54(2):66-71, 2004.
- 原田昌紀: WWWサーチエンジンの作り方, 情報処理, 40(11):1280-1283, 2000.

31