

自然言語処理入門  
「現実的なご利益: 情報検索」

東京大学 情報基盤センター  
(情報理工学系研究科、情報学府 兼担)  
中川裕志

# 情報要求

- ユーザがある問題を解決するために、現有の知識だけでは不十分であると感じている状態が情報要求(information need) な状態。
  1. 直感的要求(visceral need): 要求を認識しているが言語化できない状態
  2. 意識された要求(conscious need): あいまいな、あるいはまとまりのない表現でしか言語化できない状態
  3. 形式化された要求(formalized need): 具体的な言語表現
  4. 調整済み要求(compromised need): 問題解決に必要な情報源が同定できるくらいに具体化された要求

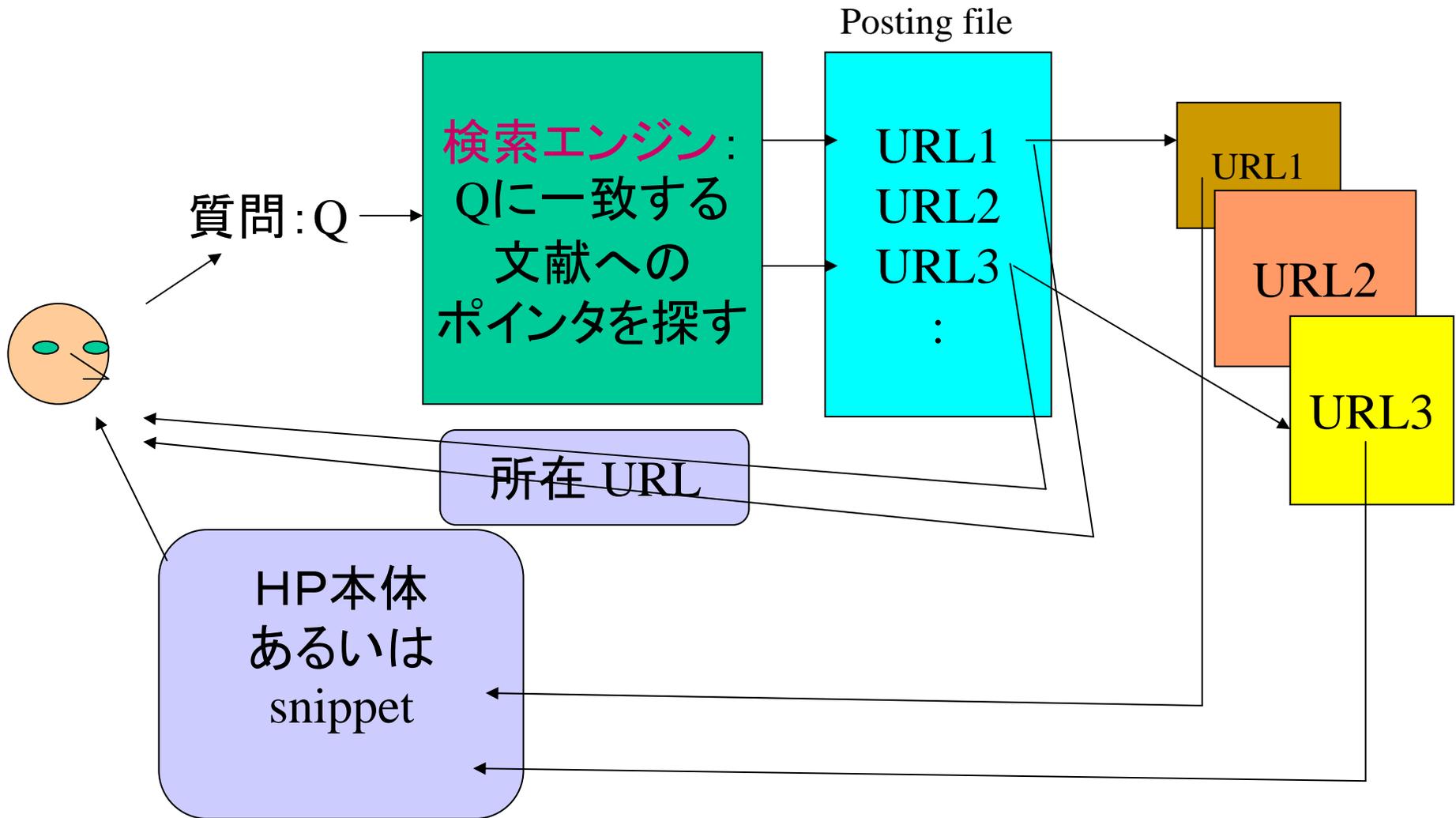
# 情報要求

- ユーザがある問題を解決するために、現有の知識だけでは不十分であると感じている状態が情報要求(information need) な状態。
- 広義の情報検索：ユーザの持つ問題を解決できる情報を見つけ出す事。問題はユーザが与える。
- 狭義の情報検索：ユーザの検索質問に適合する文書を文書集合中から見つけ出す事
- 情報の解釈をユーザが行なう= 情報検索
- ユーザがある問題を解決するために、現有の知識だけでは不十分であると感じている状態が情報要求(information need) な状態。
- 情報の解釈をシステムが行なう= 人工知能
- データベースの場合、ユーザはデータベースの構造、内容まで知った上で定型的な検索質問を行なう。つまり、文書集合を対象とする情報検索よりもさらにユーザは情報内容に通じていなければならない。実際、ユーザにとっては大変な負担であり、データベース検索の専門家が必要になったり、典型的な質問例を利用したりすることになる。

# 検索システムおよび インデクシング

検索エンジンの仕掛け

# 検索システムの構造



# 転置インデックス

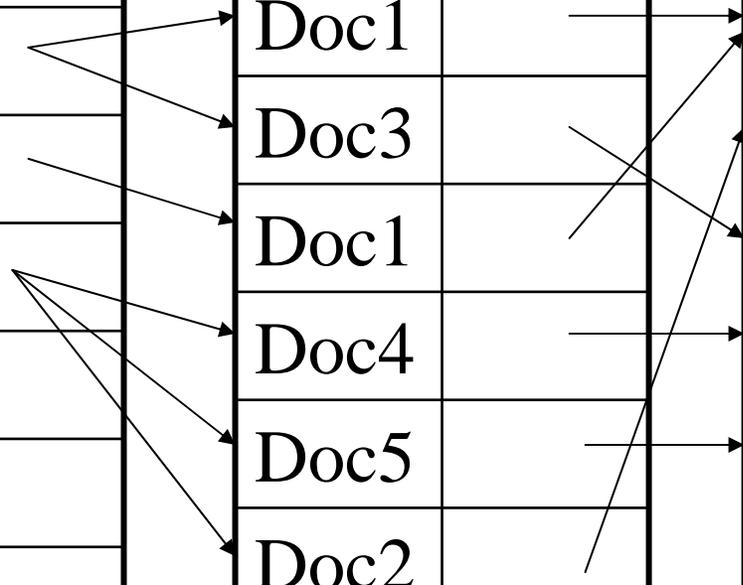
## ポスティングファイル

## Web全体

インデックスターム	ヒット件数	ポインタ
科学	2	
研究費	1	
申請	3	

URL名	ポインタ
Doc1	
Doc3	
Doc1	
Doc4	
Doc5	
Doc2	

Webページ本体
Doc1
Doc2
Doc3
Doc4
Doc5
Doc6
Doc7



# 検索システムの構造

1. 検索対象の文書各々は、それに現れる複数の語(以下「ターム(term)」と呼ぶ)によって特徴付けられる。タームによる文書の特徴付けをインデクシング(indexing) という。
2. 利用者は自分が欲しい文献を特定する記述を質問Qとして検索システムに与える。質問は、1個以上のキーワード、あるいは自然言語文で記述される
3. 検索エンジンは、Qを解釈して、それに適合する文書を検索する。適合する文書は一般的には複数存在する。したがって、検索結果は、ポスティングファイルへのポインタの集合である。ポスティングファイルは、(文書識別子、文書の格納されているアドレス) というペアからなる。
4. 文書本体は検索システム内で集中管理されている場合と、インターネット上に散在する場合がある。後者の場合、ポスティングファイルの文書格納アドレスは、URL あるいはURI である。
5. 利用者には検索結果として検索された文書あるいは文献名、URL などのアドレスが返される。多くの場合、文献名には文献のURL へアクセスできるリンクが張っており、クリックすればURL にアクセスして文献が表示できる。

# ターム と 検索

1. **統制ターム**: 予め決められたタームだけを使う。タームの質は統制により維持されるが、利用者が統制タームしか質問に使えない。
  2. **自由ターム**: 任意のタームを質問に使える。後に述べる全文検索では文書中のタームを網羅するので、自由タームになる
- ◆ 完全一致検索(ブーリアン検索)
  - ◆ 最良優先検索(ベクトル空間法など)

# ブーリアン検索

- 質問は論理式。例えば、 $Q = \text{科学} \text{or} (\text{研究費} \text{and} \text{申請})$
- 転置ファイルを検索して、ターム $t$ に対応する文書の集合 $D(t)$ が得られるとしよう。すると、質問の論理式(1)によって検索した結果は次のようになる。
- $D(\text{科学}) \text{ or } (D(\text{研究費}) \text{ and } D(\text{申請}))$
- 質問論理式を作ることが素人には難しい。
- 検索結果の文書数が多過ぎたり、あるいは0だったりすると、論理式を調整して適当な大きさの結果集合にするのだが、この調整が非常に難しい。

## 転置インデックスの高速検索

- 膨大な数のキーワードの中から質問として与えられたキーワードを高速に検索することが重要
- B木
- トライ
- PATRICIA木

# 全文検索

- メモリ中にあたかも全文が格納されているかのような状態で検索
- 本当に全文が格納されている場合＝長大な文字列マッチ処理
- 本当はなんでもかんでもインデックスされた場合＝インデックスの効率化
- 全文検索でできそうな検索質問
  - 「ターム「情報」の前後n文字以内の範囲にターム「収集」が現れる」
  - 「ターム「分析」と「統合」が同じ段落に現れる」
- インデックスの作り方
  
- 文字列を十分に長くすれば全ての半無限長文字列は異なるものになる。その上で文字列を辞書式にソートすると言語の統計で説明したKWICになる。そして、葉が反無限長文字列の開始位置をデータに持つTrieやPATRICIA木を作れば、KWICは全文検索向きのインデックスとして使える。
- 半無限長文字列の代わりに1単語を葉とするTrieやPATRICIA木にするとコンパクトになるが、葉には文書において複数回現れるその単語開始位置を全て記述しておく。これもインデックスとして使える。

# Suffix Array

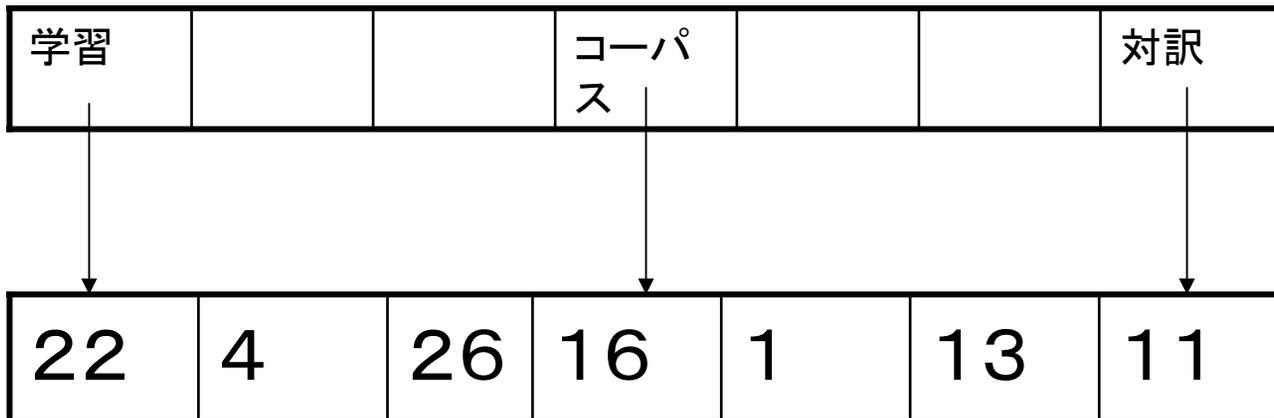
最近の機械翻訳では、対訳辞書をコーパスから学習する研究が多い

1      4                      11 13      16                      22      26

- この例文の Suffix Array は辞書式にならべた各単語の開始位置を記載した配列

22	4	26	16	1	13	11
----	---	----	----	---	----	----

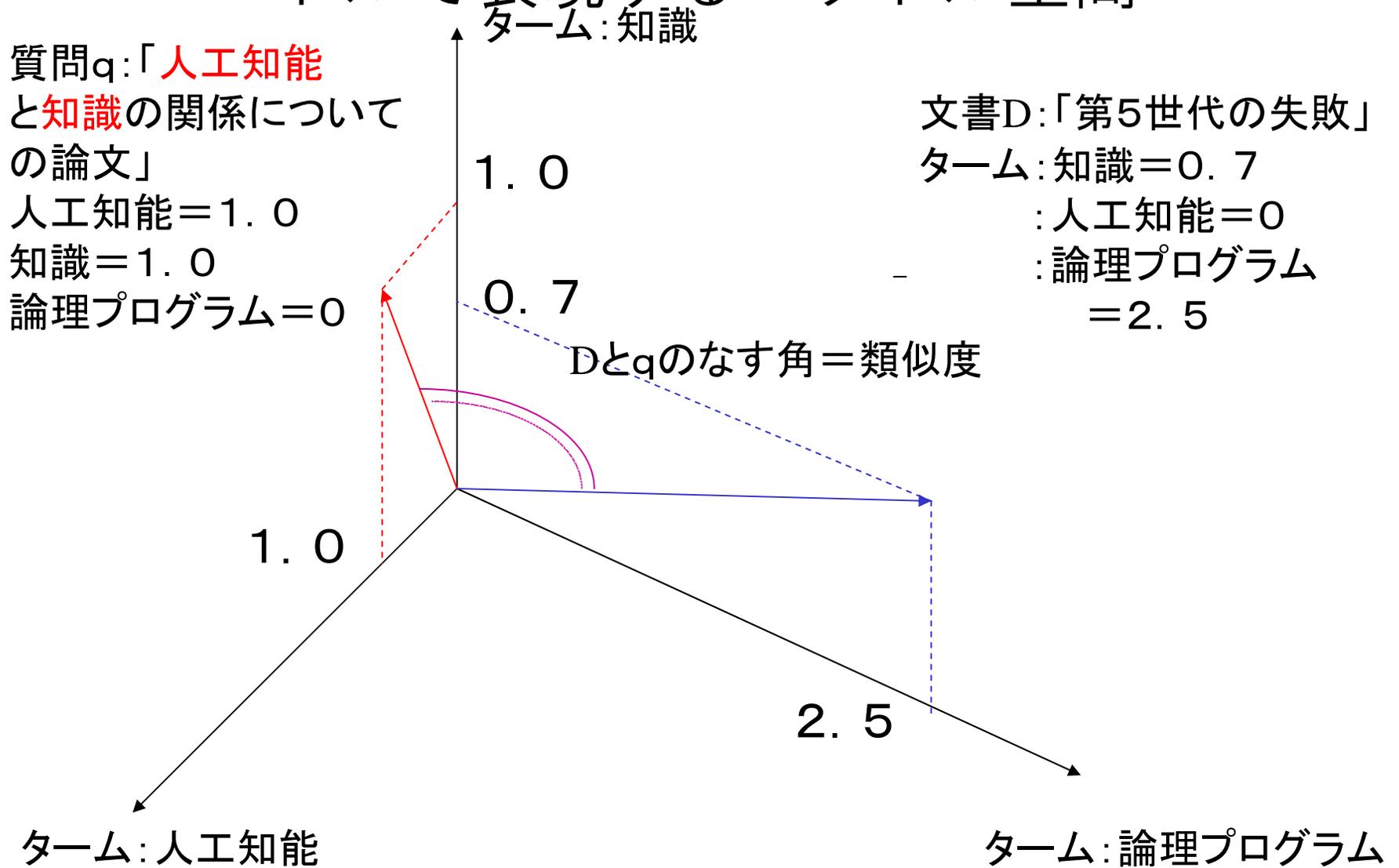
- 2分探索で高速サーチ
- Supra Index



# ベクトル空間法

- 最良優先検索
- タームの重み付けと類似度
- パッセージ検索

# 各タームを次元にし、質問と文書をベクトルで表現するベクトル空間



## タームの重み その1ターム頻度

- ターム頻度(Term Frequency: tf)  $tf_j^i = \text{freq}(i, j)$

- $\text{freq}(i; j) =$  文書Dj におけるターム $t_i$  の出現頻度。

- 変形版tf

$$tf_j^i = K + (1-K) \frac{\text{freq}(i, j)}{\max_{i, j}(\text{freq}(i, j))}$$

$$tf_j^i = \frac{\log(\text{freq}(i, j) + 1)}{\log(\text{文書 } j \text{ 中の総ターム種類数})}$$

## タームの重み その2 文書頻度

- 文書頻度 Document frequency  $df_i = Dfreq(i)$
- ただし、 $Dfreq(i)$ はターム $t_i$ が出現する文書数
- 実際はその逆数  $idf_i$  を使う
  
- 文書総数 $N$ による正規化

$$idf_i = \log \frac{N}{Dfreq(i,j)} + 1$$

## タームの重み その3 tf·idf

- 文書 $D_j$ に現れるターム $t_i$ の重み $w_{ij}$ は、 $D_j$ には数多く現れ、他の文書にはあまり現れないという性質をもつべき。つまり、文書 $D_j$ をよく特徴つけることが大切。そこで、前記のtfとidfをかけたものがよい。つまり、 tf·idf

$$w_{j}^i = \text{tf}_{j}^i \times \text{idf}_i$$

# 文書ベクトルと質問ベクトルとそれらの類似度 その1

- このようにしてターム $t_i$ の重みが決まったので、文書 $D_j$ のベクトルは、各タームを次元に割り当てた多次元空間におけるベクトルとして表現できる。つまり、

$$\mathbf{D}_j = (w_j^1, w_j^2, \dots, w_j^m)$$

- 一方質問 $q$ もターム $t_i$ を含めば1、含まなければ0という値にしてベクトルで表現できる。つまり

$$\mathbf{q} = (q_1, q_2, \dots, q_m)$$

- ただし、 $m$ は文書集合における全ての異なりターム数

## 文書ベクトルと質問ベクトルとそれらの類似度 その2

- さて、情報検索とは、質問  $q$  に対して類似度の高い文書  $D_j$  を探すことなので、類似度  $\text{sim}$  を以下に定義する。これは、ベクトル空間における  $q$  と  $D_j$  のなす角  $\theta$  が0に近いほど類似度が高いと考える方法。

$$\text{sim}(\mathbf{q}, \mathbf{D}_j) = \frac{q_1 w_j^1 + \dots + q_m w_j^m}{\sqrt{q_1^2 + \dots + q_m^2} \times \sqrt{(w_j^1)^2 + \dots + (w_j^m)^2}} = \cos \theta$$

- $\text{sim}$  の大きい順に検索結果をに並べて質問者に提示する。

# ページ間リンクを利用する重みづけ

- Webのホームページはリンクが張られている。これは図書の引用索引に似ている。
- **基本原理**：多くの良質なページからリンクされているページは、良質なページである
- この原理によるアルゴリズムとして有名なものが GoogleのPageRank algorithm と HITS algorithm

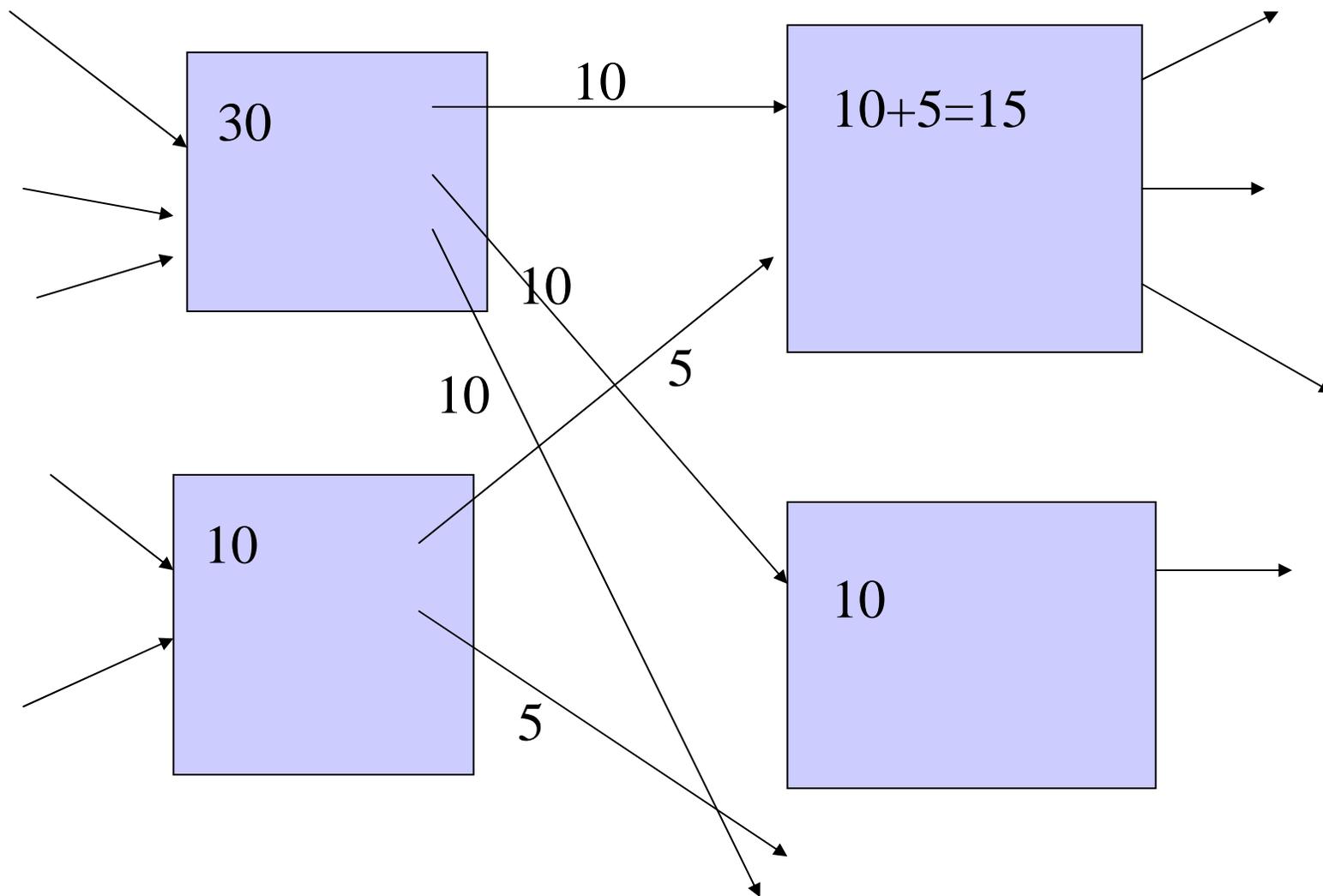
# Webページのランキング

- Webページ検索の特異性
  - 文書DBからの検索の場合、使える情報は文書中のテキスト、単語
  - $tf \times idf$ によるベクトル空間の類似度ではない重み付けはできないか？
  - Webページの検索の場合、テキスト、単語に加え、リンク情報が使える
  - リnkは(テキストや単語と異なり)、Web(=データベース)の大域的構造
  - リnkを利用する検索結果のランキングについて説明する。代表的なPageRankについて説明する。

# PageRank algorithm のアイデア

- 多くの良質なページからリンクされているページは、やはり良質なページである
- 利用者がページをリンクをたどってサーフする状況を模擬する。(良いページには沢山の利用者がいると考えればよい)
  - →
  - ページの重みが決められている
  - リンク元ページからの流入する重み総和がそのページの重み
  - リンク先へ流入する重みは、そのページの重みをリンク先の本数で割ったもの

# ページ重みの概念図



# ところが問題は

- あるページ $P_i$ の重みは $P_i$ へのリンク元のページの重みが確定しないと計算できない。
- めぐりめぐって、 $P_i$ の重みが確定しないと $P_i$ へのリンク元のページの重みが確定しないと計算できない。
- →
- 各ページの重みを要素とし、全ページ数の次元を持つベクトル  $R$  を考え、このベクトルを数学的に求めることにする

# ページ重みの計算法

- 接続行列  $A=(a_{ij})$  で表現

- if ページ  $i$  から ページ  $j$  にリンク then  $a_{ij}=1/N_i$

- ( $N_i$  は  $i$  から できるリンク数)

- otherwise  $a_{ij}=0$

- ページのベクトル  $R$  は以下の式を満足する

- $R=cAR$  つまり

$$\begin{bmatrix} r_1 \\ \vdots \\ r_n \end{bmatrix} = c \begin{bmatrix} a_{11} & \cdots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{n1} & \cdots & a_{nn} \end{bmatrix} \begin{bmatrix} r_1 \\ \vdots \\ r_n \end{bmatrix}$$

# しかし、その1

- 2つのページの間だけでループしている状態に他からのリンクがあり、そこに重みがどんどん吸収されてしまう場合
- ページがいくつかのグループにまとめられ、グループ間にはリンクがない
- このような場合には、全ページにわたっての重みが均等に行き渡らない。
- →ある割合(15%)で、ランダムにページを探したと仮定
- →Aに全要素が一定値の行列を加算して計算
- ◆  $R = (cA + (1-c)E)R$ 
  - ◆ Eは全要素が1の行列、 $c=0.85$ としている

## しかし、その2

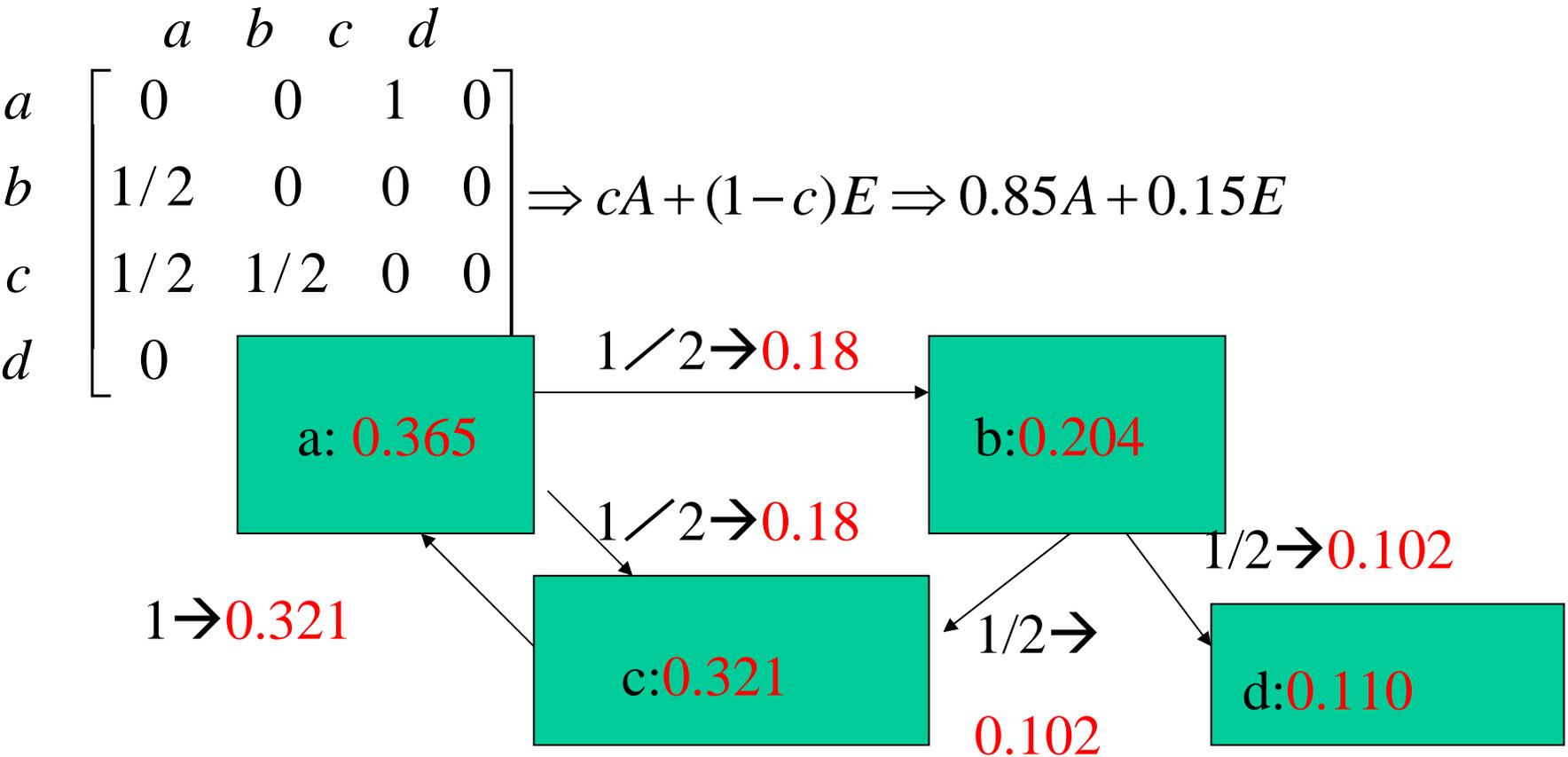
- 流出すれだけで、流入する全くリンクのないページとは
- 勝手に他人にリンクを張っているだけの権威のないページ → 無視する
- 機械的に張ったリンクの問題
  - PageRankのアイデアの基本は、人間がWebページを評価しリンクを張った状況、すなわち人間の判断を利用しようとしているので、機械が張ったリンクは悪影響を与える

# ページベクトルの計算

- ◆  $R = (cA + (1-c)E)R$
- ◆ Eは全要素が1の行列、 $c=0.85$ としてみた
- ◆ これはベクトルの固有値問題でRは固有ベクトル。これを以下の繰り返し計算で解く。
  1.  $R(0) \leftarrow$  適当な初期値
  2.  $R(i+1) \leftarrow (cA + (1-c)E) R(i)$
  3.  $D \leftarrow \|R(i)\| - \|R(i+1)\|$
  4.  $R(i+1) \leftarrow R(i+1) + dE$
  5.  $\delta \leftarrow \|R(i+1) - R(i)\|$
  6. if  $\delta > \varepsilon$  then goto 2
- ◆ なお  $\|x\|$  はベクトルの全要素2乗の和の $\sqrt{\quad}$

# 計算例

- 簡単な例で繰り返し計算すると( $\varepsilon = 0.001$ )
- Rは  $\sum_{i=1}^4 r_i = 1$  と確率になるように正規化



# 考察

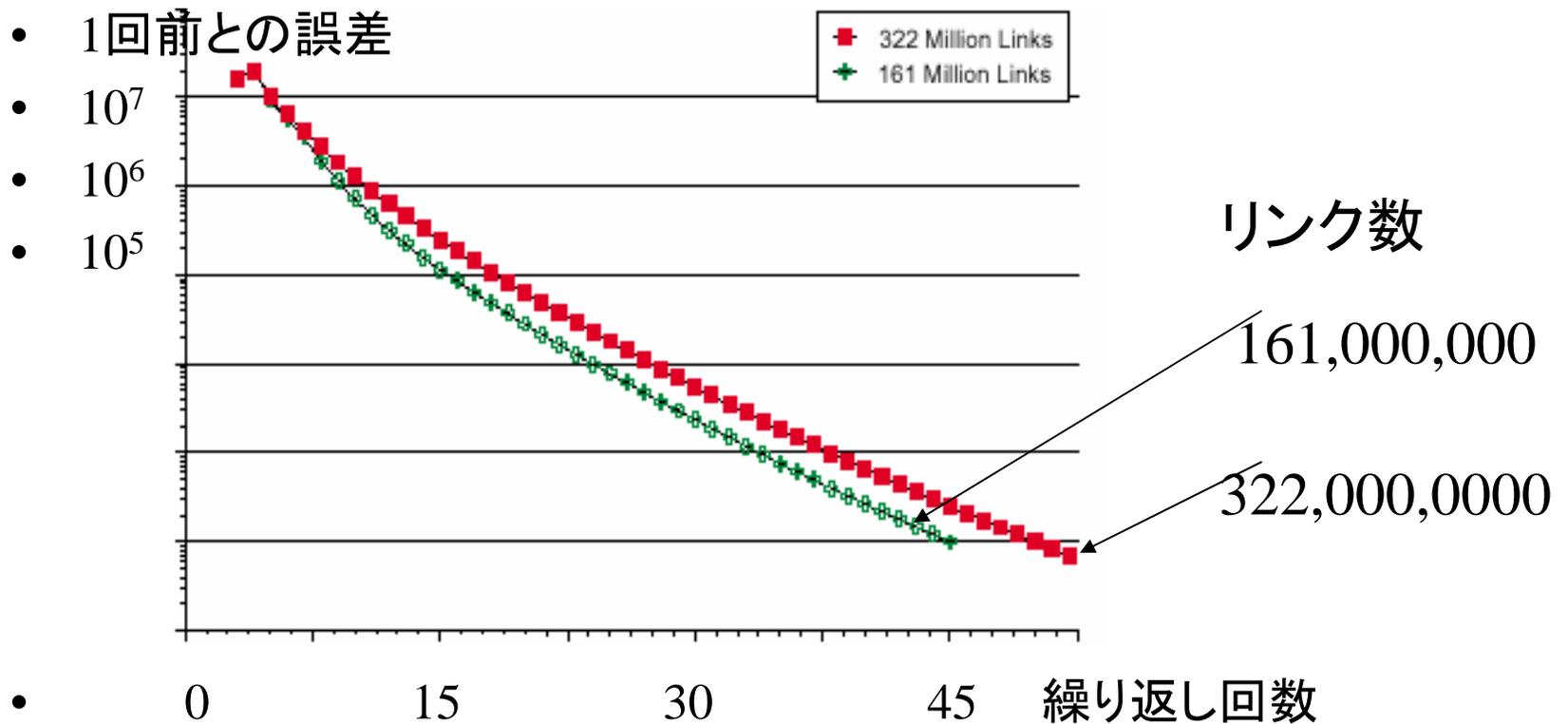
- I. a,b,c,d ページともほぼ流れ込むリンクの重みに近い
- II. リンク重みとページ重みのズレは、0.15の確率でランダムにページを移動する分
- III. 良いページcからリンクされているaは重みが大きい
- IV. 多くのページからリンクされているcも重みが大きい
- V. あまりたいしたことのないページbだけからリンクされているdは重みが小さい
- VI. 以上、ほぼ我々の意図したページの重みになっている。

# 数値計算の問題

- ✓ これまで述べてきたページ重みの繰り返し計算はうまく収束するのか、十分早いのか？
- ✓ Googleの10億ページでは100億のリンクがあるので、まっとうに線形代数の計算をしていたのでは無理。
  - ✓ 数値計算法ではかなり工夫をしているはず。企業秘密らしく公開されていない
  - ✓ Googleのマシナールームの様子からみると、
    - ✓ PCクラスタによる並列計算
  - ✓ 疎な行列を適当に分解すればまとまったページ群を独立に計算できるかも

# 収束の問題(以下の論文の数値例グラフ)

Lawrence Page, Sergey Brin, Rajeev Motwani, Terry Winograd,  
'The PageRank Citation Ranking: Bringing Order to the Web', 1998



# HITS algorithmのアイデア

- ◆ 探し当てたいWebページは必ずしも検索のキーワードを含まない
  - ◆ TOYOTAやHONDAのページに自動車製造業なんてキーワードは見当たらない
  - ◆ 一方、yahooなどの権威あるページはこういうキーワードをむ
  - ◆ キーワードではなくWebのリンク構造を利用 (PageRankと同じ意図)
- ◆ HITSは、ある話題に対するauthoritiesと多数のauthoritiesにリンクを張る hub との関係を使う方法

# Focused Subgraph - 1

- ▶ ある話題に関するページ集合:  $S$  が対象
- ▶  $S$  の持つべき性質は
  - ▶ 比較的小さい
  - ▶ 関連あるページを十分に含む
  - ▶ ほとんど全ての (あるいは多数の) authority page を含む
- ▶ トピック:  $Q$  を表わすキーワードで検索エンジンから拾ってきたページ集合では、上記の条件は満たさない。

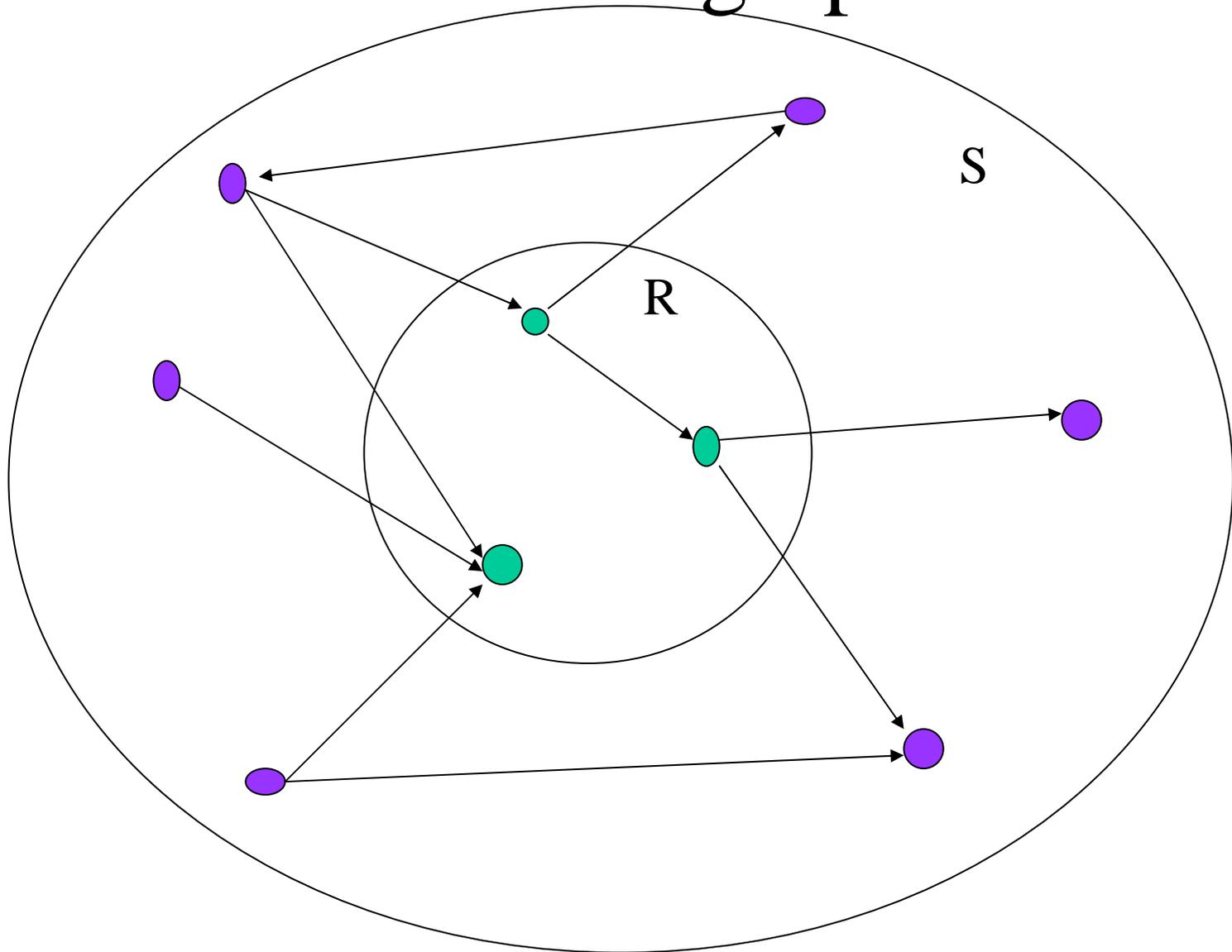
# Focused Subgraph - 2

- ▶ 例えば、検索結果の上位20ページを集めると
- ▶ javaというキーワードで検索すると、別のサイトのページ間には15links
- ▶ censorshipというキーワードで検索すると、別のサイトのページ間には28links
- ▶ 最大link数 =  $200 \times 199 = 39800$ には遠く及ばない
- ▶ そこで、

# Focused Subgraph — 3

- 比較的少数のタネ集合:Rからはじめる。
- Rは質問:Qの検索結果の上位  $t$  ページ
- ①  $S=R$
- ② Rに属するページからリンクしている(out going link)ページを  $S$  に追加
- ③ Rに属するページにリンクされている(in coming link)ページのうち  $d$  ページをランダムに選んで  $S$  に追加
- Kleinbergの論文では  $t=200, d=50$ で $S$ は 1000 ~5000ページ

# Focused Subgraph — 4



# Hub と Authority の計算

- ◆ S中には、トピック:Qに関係がないのに多くのリンクが入ってくるページがあり、
- ◆ 重要なページでないのに、やたらと他のページにリンクを張りまくっているページがある
- ◆ 真の authority と hub は相互に強化する
  - ◆ Good hub pages points to many good authority pages, good authority pages are pointed to by many hubs.
  - ◆ 同じdomain名の中でのリンクの張り合い無視（目次や自己参照などだから）

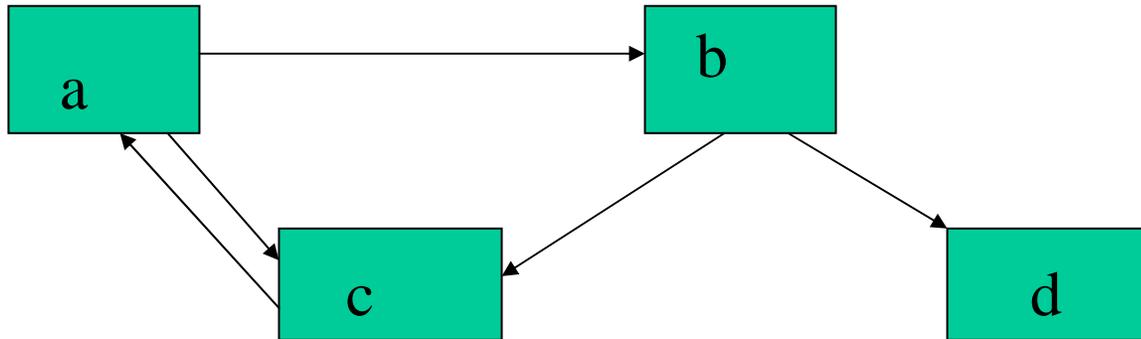
# weights

- Authority weight of page:  $\mathbf{p} = \mathbf{y}^{\langle p \rangle}$
- Hub weight of page:  $\mathbf{p} = \mathbf{x}^{\langle p \rangle}$
- Link set:  $E = \{(p, q) \mid \text{page : } p \text{ points to } \text{page : } q\}$
- Normalization  $\sum_{p \in S} (\mathbf{x}^{\langle p \rangle})^2 = 1, \quad \sum_{p \in S} (\mathbf{y}^{\langle p \rangle})^2 = 1$

# iteration

- ① :  $x_0 = (1,1,\dots,1) \in \mathbf{R}^n$ ,  $y_0 = (1,1,\dots,1) \in \mathbf{R}^n$
- ② : 以下③から⑥をk回繰り返す(k=20くらいで十分)
- ③ : *for every*  $p \in S$   $x^{<p>} = \sum_{q:(q,p) \in E} y^{<q>}$
- ④ : *for every*  $p \in S$   $y^{<p>} = \sum_{q:(q,p) \in E} x^{<q>}$
- ⑤ :  $\{x^{<p>}\} \leftarrow \textit{normalized}$   $\{x^{<p>}\}$
- ⑥ :  $\{y^{<p>}\} \leftarrow \textit{normalized}$   $\{y^{<p>}\}$
- ⑦ : 結果は  
 $\textit{authority value} = x^{<p>}$   
 $\textit{hub value} = y^{<p>}$

# 計算例



● Authority:  $a=3.46 \times 10^{-6}$ ,  $b=0.408$

$c=0.816$ ,  $d=0.408$

● Hub:  $a=0.707$ ,  $b=0.707$ ,

$c=5.98 \times 10^{-6}$ ,  $d=0$

$$A = \begin{matrix} & \begin{matrix} a & b & c & d \end{matrix} \\ \begin{matrix} a \\ b \\ c \\ d \end{matrix} & \begin{bmatrix} 0 & 0 & 1 & 0 \\ 1 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{bmatrix} \end{matrix}$$

# 数学的再考－1

◆ 結局はページ間の接続行列に関連する固有値問題

◆ 接続行列

$$A = \begin{bmatrix} a_{11} & \cdots & a_{1n} \\ \vdots & \ddots & \\ a_{n1} & \cdots & a_{nn} \end{bmatrix}$$

$$a_{ij} = \begin{cases} 1 & \text{if page } i \text{ points to page } j \\ 0 & \text{otherwise} \end{cases}$$

◆ 繰り返し計算は

$$\mathbf{x}_i = \begin{pmatrix} x_i^{<1>} \\ \vdots \\ x_i^{<n>} \end{pmatrix} \quad \mathbf{y}_i = \begin{pmatrix} y_i^{<1>} \\ \vdots \\ y_i^{<n>} \end{pmatrix}$$

$$\mathbf{x}_{i+1} \leftarrow A^t \mathbf{y}_i \quad \text{and} \quad \mathbf{y}_{i+1} \leftarrow A \mathbf{x}_{i+1}$$

## 数学的再考－2

◆ したがって

$$\mathbf{x}_k = (\mathbf{A}^t \mathbf{A})^{k-1} \mathbf{A}^t \mathbf{z}, \quad \mathbf{y}_k = (\mathbf{A} \mathbf{A}^t)^k \mathbf{z}, \quad \mathbf{z} = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}$$

◆ ゆえに

$\mathbf{A}^t \mathbf{A}$  は  $n \times n$  の対称行列だから、Golub and Van Loan の線形代数の結果により、

$\mathbf{x}^\infty = \mathbf{A}^t \mathbf{A}$  の最大固有値に対応する固有ベクトル

$\mathbf{y}^\infty = \mathbf{A} \mathbf{A}^t$  の最大固有値に対応する固有ベクトル

# 实例

Jon M. KlienberglのAuthorities Sources in a Hyperlinked Environment, JACM 46-5 による

- (java) Authorities
- 0.328: <http://www.gamelan.com>
- 0.251: <http://java.sun.com/>
- 0.190: <http://www.digitalfocus.com/digitalfocus/faq/howdoi.html>
- 0.190: <http://lightyear.ncsa.unic.edu/~srp/java/javabooks.html>

# パッセージ検索

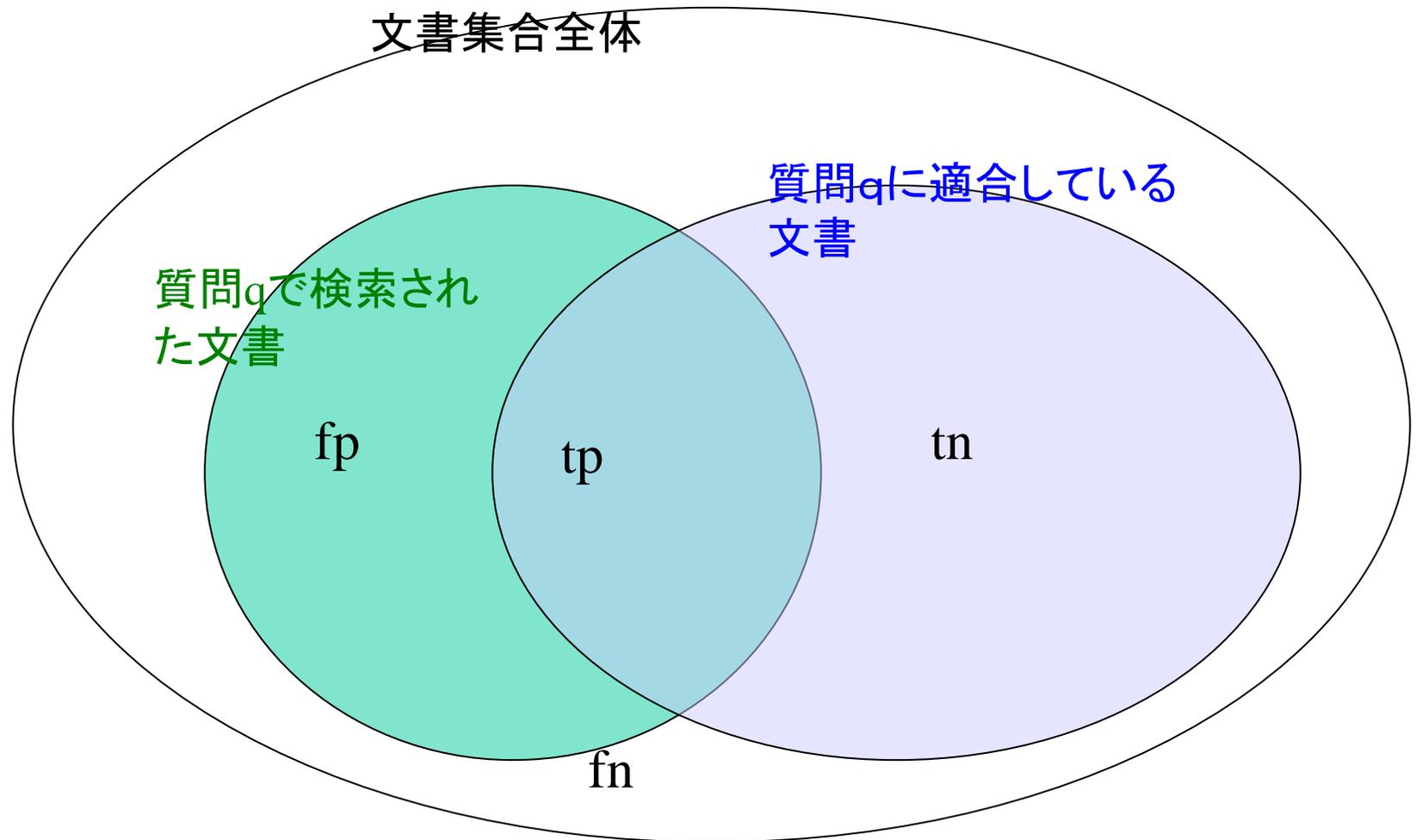
- 文書の内容を特徴付けるのは文書全体よりはむしろ特定の部分
- ベクトル空間モデルを文書ではなく、文書の小さな部分、例えば段落、に適用。この小さな部分をパッセージという。つまり、文書Dの代わりにパッセージPkを使って、パッセージ重み $w_k^l$ を計算し、ベクトル空間法を適用
- パッセージの候補としては、
  - 1 固定長に分割したテキストの部分
  - 2 形式段落
  - 3 形式的な節、章

# 情報検索システムの評価法

- 評価法
- 再現率、適合率
- F値など

## 一般的な検索結果の状態

- 質問 $q$ で結果の文書集合が得られた。しかし、結果の中には間違いもあるし、得られなかった文書の中にも正解がありうる。



# 検索エンジンの性能評価

- 再現率

$$recall = \frac{tp}{tp + tn}$$

- 適合率あるいは精度

$$precision = \frac{tp}{tp + fp}$$

- フォールアウト

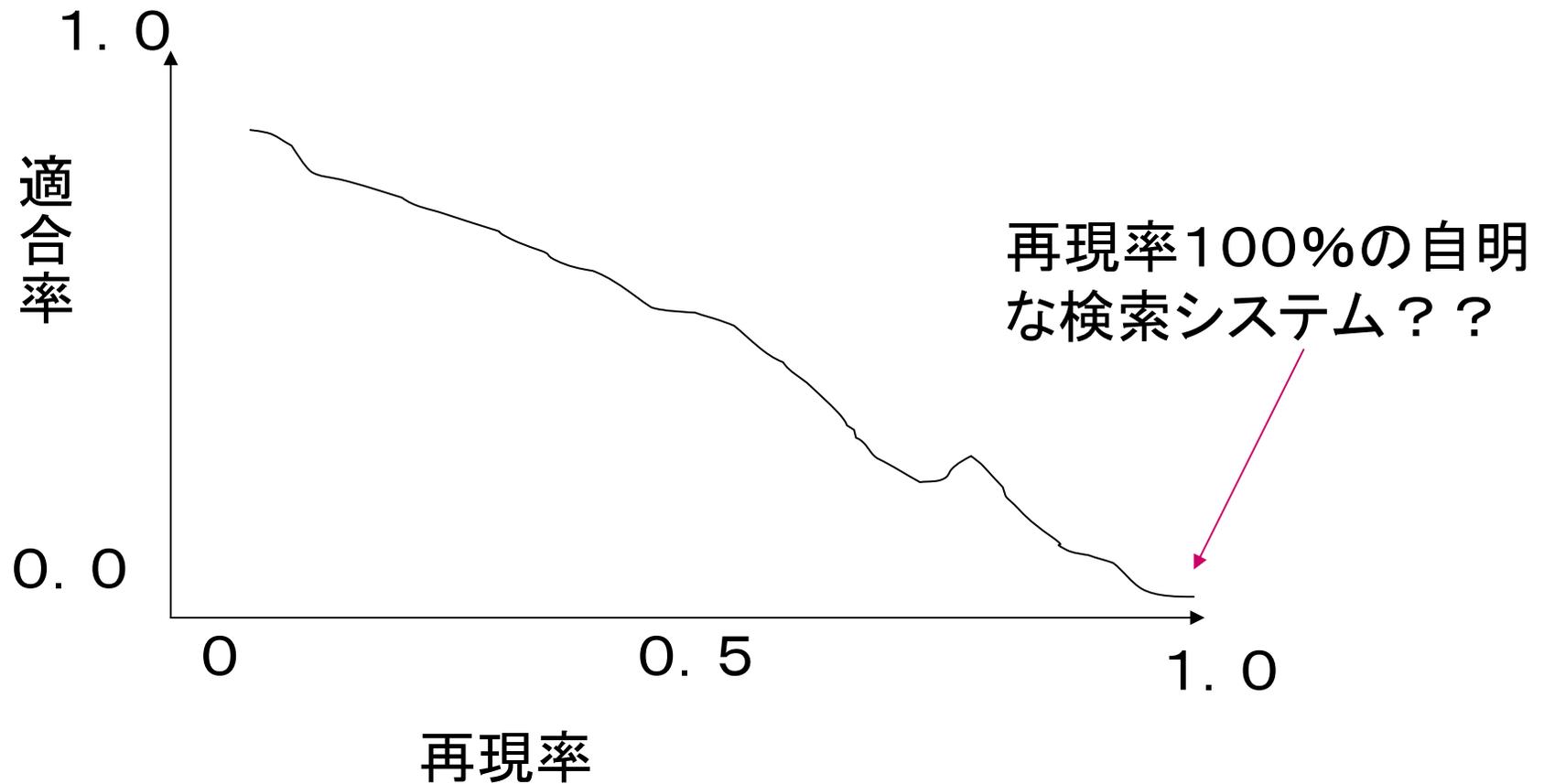
$$fallout = \frac{tn}{tn + fp}$$

- 一般性

$$generarity = \frac{tp}{tp + tn + fp + fn}$$

# 再現率 vs 適合率

- よく使う評価の表現法



## 再現率 vs 適合率に関連した尺度

- Break even point 再現率と適合率が一致する点
- 11点平均適合率 再現率=0.0, 0.1, 0.2, ..... 0.9, 1.0 の11点における適合率の平均値
- F値 ただし、bは適合率が再現率よりどれだけ重視されているかを示すパラメター b=1がよく使われる。

$$F = \frac{(1 + b^2) \times P \times R}{b_2 \times P + R}$$

$$F = \frac{2PR}{P + R} = \frac{2}{\frac{1}{P} + \frac{1}{R}}$$

# 順位つき検索結果の評価

- ブーリアン検索では検索結果は全て同等
- ベクトル空間法やPageRank、HITSでは検索結果が質問に適合した順番に並ぶ。(表示も適合順)
- この場合の評価法について

# Recall , Precision

- 質問 $q$ に適合する結果(以下、正解、という)の数:  $|D_q|$
- 検索エンジンの順位つけられた結果:  
( $d_1 \dots \dots d_n$ )
- $d_i$ が質問 $q$ への正解なら  $r_i=1$ 、 そうでなければ  $r_i=0$   
とする。すると、
- 第 $k$ 順位まで拾ったときの

$$\text{Recall}(k) = \frac{1}{|D_q|} \sum_{1 \leq i \leq k} r_i$$

$$\text{Precision}(k) = \frac{1}{k} \sum_{1 \leq i \leq k} r_i$$

# 平均適合率 : average precision

$$\text{AveragePrecision} = \frac{1}{|D_q|} \sum_{1 \leq k \leq N} r_k \times \text{precision}(k)$$

ただし、 $N$ は正解が最後に現れた順位

- 例:

$$\begin{aligned} & \text{AvPrec} \\ &= \frac{1}{2} \left( \frac{1}{1} + \frac{2}{4} \right) \\ &= 0.75 \end{aligned}$$

順位	正解か
1	○
2	
3	
4	○
5	
6	

# 平均逆順位 : Mean Reciprocal Rank(MRR)

$$MRR = \frac{\sum_{n=1}^N (if \quad n\text{-位の結果が正解} \quad \frac{1}{n} \quad else \quad 0)}{\sum_{n=1}^N \frac{1}{n}}$$

ただし、 $N$ は正解数

もし、正解がひとつも 現れなければ  $MRR=0$

- 例

$$MRR = \left( \frac{1}{1} + \frac{1}{4} \right) / \left( \frac{1}{1} + \frac{1}{2} \right) = 0.833$$

順位	正解か
1	○
2	
3	
4	○
5	
6	

# テストコレクション

- (a) 文書集合、(b) 多数の質問、(c)各質問に対する適合文書の集合、を組にしたデータベースをテストコレクションと呼び、情報検索システムの性能評価において必須の資源である
- 正解集合を作ることは大規模テストコレクションでは困難
- Pooling method:、同一の文書集合に対して、多数の検索エンジンで同じ質問を出し、上位N個の検索結果を全て集める。Nの値として、100程度が多い。この結果に対してのみその適合性を人手で判断し、それを文書集合全体における適合した文書とする