

協調検索エンジンにおける 分散型スコアリング

上原、山本、佐藤、西田、森
東洋大学工学部情報工学科

目次

- 背景及び目的
- 協調検索エンジンの概要
- 分散型スコアリング
- 実装
- 評価
- まとめ

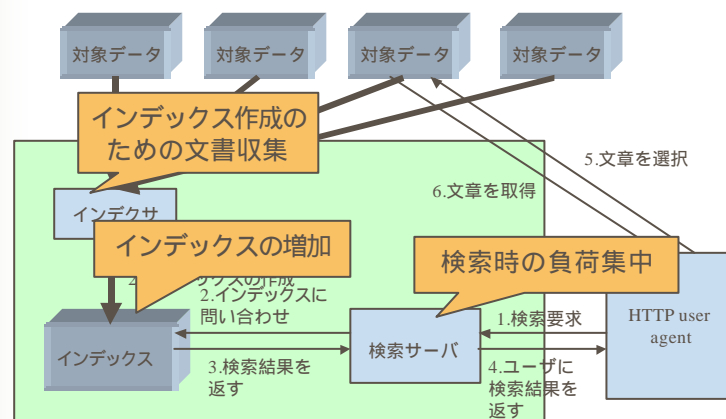
背景および目的

- 組織内での情報検索
 - 更新期間の短縮 分散型検索へ
- 協調検索エンジンの提案
 - 検索エンジンによりスコアが異なる
 - 同種の検索エンジンでも $tf*idf$ 計算は正しく行えない
- 分散型 $tf*idf$ スコアリングが必要

1999/12/23

3

集中型全文検索システム



1999/12/23

4

協調検索エンジンの提案

■ 概要

- 複数の検索エンジンを協調させ、1つの検索エンジンのように見せる

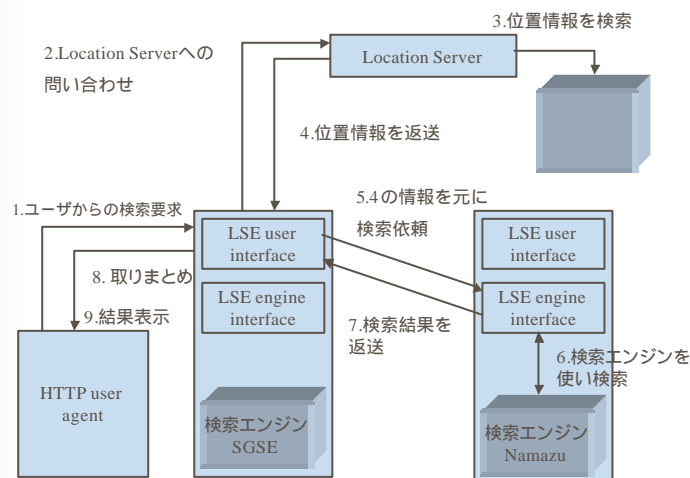
■ 目的

- 組織内情報検索の統一
- 既存検索エンジンとの互換性
- 更新期間の短縮
- 集中型検索エンジンに劣らない性能

1999/12/23

5

協調サーチエンジン



1999/12/23

6

協調検索エンジンの問題点

- 基本検索エンジン(Namazu、SGSE)によりスコアが異なる
 - スコアの共通化が必要
- 同種のエンジンでもtf*idfスコアリングが正しく行えない



- 分散型tf*idfスコアリングの実現

1999/12/23

7

tf・idf法

- $s(d,k) = K \times tf(d,k) \times idf(k)$
 - $s(d,k)$: キーワードkに対する文書dのスコア
 - K: タグなどによる重み付け
 - $tf(d,k)$: 文書dにおけるキーワードkの出現頻度
 - $idf(k) = \log(N/n)$
 - N: 総文書数
 - n: キーワードkを含む文書数

キーワードの出現頻度だけでなく、珍しさも考慮したスコアの計算法

1999/12/23

8

分散tf・idf法

■ $s(d,k) = K \times tf(d,k) \times idf(k)$

- $s(d,k)$: キーワードkに対する文書dのスコア
- K: タグなどによる重み付け
- $tf(d,k)$: 文書dにおけるキーワードkの出現頻度
- $idf(k) = \log (N_i / n_i)$
 - N_i : サイトiの総文書数
 - n_i : サイトiのキーワードkを含む文書数

N_i, n_i を集計する必要がある

1999/12/23

9

NamazuとSGSEにおける 重み付けの相違

タグの種類	Namazu の重み	SGSE の重み
Keywords	32	100
description	32	10
title	16	2
H1~6	8~3	1
A	4	1
STRONG,EM,KBD,SAMP,VAR,CO DE,CITE,ABBR,ACRONYM,	2	1
その他	1	1

1999/12/23

10

NamazuとSGSEにおける 論理型検索におけるスコア計算

■ Namazu

- $s(d,k) = K * tf(d,k) * idf(k)$
- $s(d,A \text{ and } B) = \min(s(d,A), s(d,B))$
- $s(d,A \text{ or } B) = \max(s(d,A), s(d,B))$
- $s(d,A \text{ not } B) = s(d,A)$

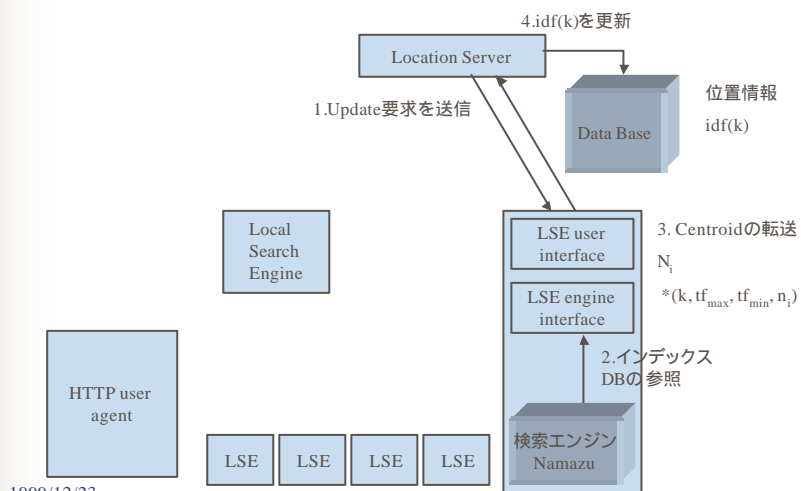
■ SGSE

- $s(d,k) = K * tf(d,k)$
- $s(d,A \text{ and } B) = s(d,A) * s(d,B)$
- $s(d,A \text{ or } B) = s(d,A) + s(d,B)$
- $s(d,A \text{ not } B) = s(d,A)$

1999/12/23

11

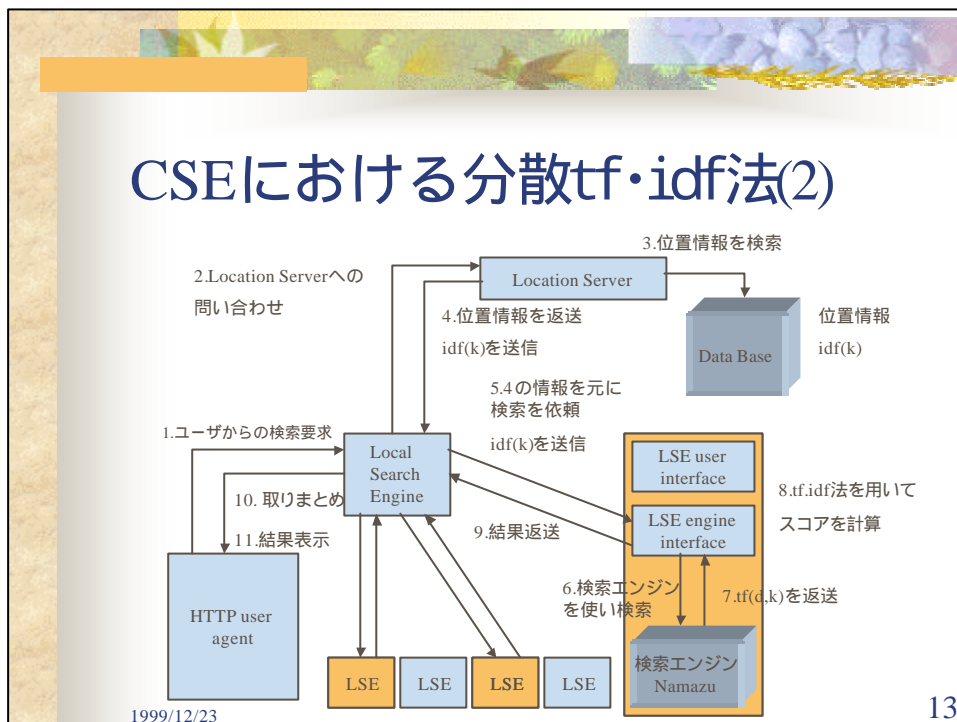
CSEにおける分散tf・idf法(1)



1999/12/23

12

CSEにおける分散tf·idf法(2)

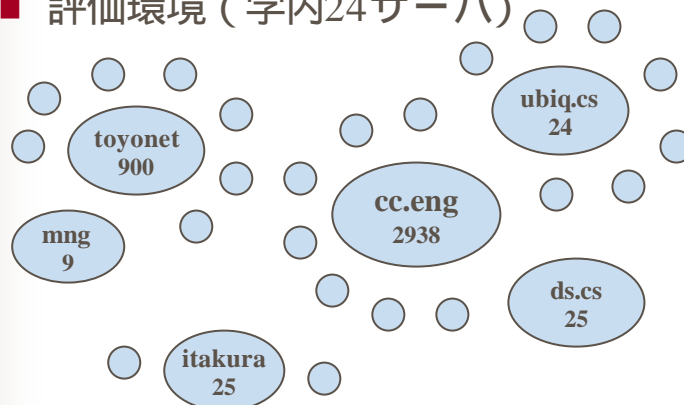


実装

- 方針
 - スコアリングの基準をNamazuに合わせる
 - インデックス情報のみを変換し、検索エンジン自体は変更しない(互換性・相互運用性)
- 方法
 - メタエンジン(LSE)が検索論理式を単独キー検索に変換し、論理型検索をエミュレーションする

実際の環境における評価(1)

■ 評価環境 (学内24サーバ)



1999/12/23

15

実際の環境における評価(2)

	応答時間 (秒)
CSE	1.267
Namazu	0.046
SGSE	0.448

単独サイトにおける単一
キーワードによる最短検
索時間の比較

		応答時間 (秒)
AND検索	CSE	1.5
	Namazu	5.6
OR検索	CSE	5
	Namazu	5.6

CSEの並列検
索とNamazuを
用いた人手に
よる逐次検索
との比較

1999/12/23

16

まとめ

- $tf*idf$ に基づくスコアを分散計算
- 既存検索エンジンとの互換性・相互運用性を確保
- 通信遅延は大きいですが、サイト集合の絞り込みにより実用的な性能を達成
- インデックスの更新時間は大幅に短縮

1999/12/23

17

今後の課題

- 更新時間の更なる短縮
 - 高速なインデクサーの開発
- キャッシュサーバ
 - 次の10項目を素早く返す
 - 規模の拡大(組織内から組織間へ)
- 分散型収集ロボット

1999/12/23

18