



文書クラスタリングの技法

7. Leader-follower アルゴリズムの応用
8. Crouch のアルゴリズム
9. CBC アルゴリズム



Leader-followerアルゴリズム

- **特徴**
- クラスタ個数の代わりに文書とクラスタの類似度の閾値を設定しておく
- 分類対象を処理していく中で閾値に基づいて自動的にクラスタが形成されていく



単純なleader-follower法①

(1) 閾値 θ_s を設定して最初の文書 d_1 を読み、最初のクラスタとする。

→ クラスタの個数 L は 1 になる。

→ i を初期化しておく。

(2) $i \leftarrow i+1$ とし文書 d_i を読む。すべての文書を読み終わっていれば終了。そうでない場合は(3)に進む。



単純なleader-follower法②

(3) 文書ベクトル d_i とその時点での全てのクラスタベクトル c_k との類似度 $s(d_i, c_k)$ を計算する。

もし閾値を越えれば d_i をクラスタ c_k に加え、ベクトル c_k を更新する。

もし文書がどのクラスタにも加わらなければ新しいクラスタとする。

$\rightarrow L \leftarrow L + 1, c_L \leftarrow d_i$

この後(2)に戻る。



単純なleader-follower法③

- **問題点**
- θ sの設定が課題となっている
- L が大きくなると計算量が増加する
- クラスタベクトルの格納に多くのメインメモリを消費する



Crouchのアルゴリズム ①

(1) クラスタの設定

(2) クラスタへの文書の割り当て

という二段階からなる。

段階(1)

- ・ベクトル ck を更新する際に各語の変動係数を計算してそれに基づき各クラスタを表現する語の集合(core terms)を選別する。



Crouchのアルゴリズム ②

- ・core termsの重さの平均に基づいてベクトルckを更新する。

段階(2)

- ・クラスタの重複が許され、一件の文書が複数のクラスタに属することが可能である。

利点

- ・クラスタベクトルを格納する領域が少ない
- ・類似度sの計算時間が短い



CBCアルゴリズム①

段階Ⅰ 文書ごとに類似度の高い20文書を求める

段階Ⅱ 段階1で計算された類似度がcommittee
リストを求める。

段階Ⅲ 各文書をもっとも類似したcommitteeに割
り当てる。

* ここでcommitteeとは各クラスタの核となるいわ
ば「代表者」であり、段階2で求められる。



CBCアルゴリズム②

段階Ⅱ

1. 2つの閾値 θ_1, θ_2 を決める。3種のリスト L_c, L_k, L_r を用意して L_k を committee リストとする。
また、未処理の文書をリスト E に格納する。
2. リスト L_c を空にし、 E に含まれる文書で次の処理を行う。



CBCアルゴリズム③

2-1 その文書と類似した20文書に対して群平均クラスタリングを適用する。

2-2 その結果生成されたクラスタCごとに得点 $|C| \times \text{avgsim}(C)$ を計算する。

- ・ $|C|$ はそのクラスタに含まれる文書数
- ・ avgsim はC中の文書の組ごとの類似度の平均



CBCアルゴリズム④

- 2-3 得点の最も高いクラスタをリストLcに追加登録する。
- 3 リストLcを得点の降順に並べ替える
- 4 リストLkを空にする(初期化)
- 5 リストLc中の上位から順に書くCについて次の処理を行う。



CBCアルゴリズム⑤

- 5-1 クラスタCの重心とその時点でリストLkに含まれている全てのcommitteeの重心との類似度を計算する。
- 5-2 それらの類似度が全て閾値 θ_1 を下回っている場合、CをcommitteeのリストLkに追加する。
- 6 リストLkが空ならば処理を終了する。



CBCアルゴリズム⑥

- 7 再びEに含まれる文書ごとに以下の処理を行う。
 - 7-1 もしその文書とすべてのcommitteeとの類似度が閾値 θ_2 を下回るならばその文書をリストLrに追加する。
- 8 もしリストLrが空ならば処理を終了する。
そうでなければLkの内容を保存し、LrをEとして2の処理に戻る。



CBCアルゴリズム⑦

出力結果

段階8で保存されたリストLkの和集合

段階Ⅲ 和集合に含まれる各committeeの重心に対する各文書の類似度を計算し、類似度の最も高いcommitteeに割り当てる。



CBCアルゴリズム⑧

利点

各クラスタの核をcommitteeとして求めることによって「周辺の」な文書がクラスタの重心の計算に不当な影響を及ぼすことを防いでいる。