



# 文書クラスタリングの技法

---

C.次元縮約法に基づくクラスタリング

1.PDDPアルゴリズム

2.LSIの応用とIRRアルゴリズム



# PDDPアルゴリズムとは？

---

- 主成分分析を利用して、文書集合を逐次的に分割していくアルゴリズムである。
- 2分木としてクラスタが階層的に構成されることになるのでPDDPは分割型の階層的アルゴリズムの一種であるといえる。



# アルゴリズムの手順(1)

---

1. 行列 $W$ の各要素を以下で定義し、各文書の長さを1に標準化する。

$$W_{ij} = \frac{x_{ij}}{\sqrt{\sum_{j=1}^M x_{ij}^2}}$$

2.  $W$ に対する転置行列の共分散行列

$$S = (W^T - ue^T)^T (W^T - ue^T)$$

を対角化する行列  $G^T S G = \Lambda$  を求める。

\* ここで $u$ は $M$ 次元平均ベクトルを意味する。



## アルゴリズムの手順(2)

---

3. 行列G中のN次元列ベクトルを $g_1$ と表記する。  
これは語 × 文書の重み行列に対して主成分分析を実行し、「寄与率」の最も大きな第1主成分を求めたことになる。
4. 固有ベクトル $g_1$ の第i番目の要素は、文書ベクトル $d_i$ を第1主成分へ射影した値なので、この値が正であるグループと負である2つのグループに文書集合を分割できる。



## アルゴリズムの手順(3)

---

5. 分割された2つの文書集合のそれぞれに対して同様な処理を反復的に繰り返す。

以上の手順で文書集合に関する2分木を構成することができる。



# LSIとは？

---

情報検索の向上のために行列 $W$ に対して特異値分解(SVD)を施し、文書空間の次元を圧縮する手法である。

語×文書の行列  $W^T$  に対する特異値分解は

$$W^T = UQV^T$$

とかける。



# 変数解説

変数名	意味
U	$M \times N$ の直交行列
Q	$N \times N$ の対角行列
V	$N \times N$ の直交行列

\* Wのランクを $r$ とするとQにおける $N-r$ 個の対角要素は0である。



# LSIの問題点

---

LSIを文書クラスタリングに応用する際は、規模の小さなクラスタがうまく抽出されない。

理由は本来LSIは全体的に影響力の小さい語を排除することによって検索機能の向上を図っている。しかしこれは同時に他の文書とあまり類似しないような文書が除かれる傾向にあるから。



# IRRアルゴリズムの概念

---

- LSIの問題点を補うために考案された。
  - 元の行列から次元を逐次的に抽出することにより、それより前の段階で取り出された次元に関連しない部分を説明するような次元の選択を可能にしている。
  - 前段階で取り出された次元を差し引いた剰余を計算し、そこから次の次元を抽出する。
- \* 行列をここでは $R$ と表記し書く列を $r_i$ と書く。



# アルゴリズム(1)

---

- (1)  $R = W^T$  とする。また  $q$  をある値に定め、 $k \leftarrow 1$  とする。さらに前もってクラスタ個数  $L$  を決めておく。
- (2)  $R_s$  を求める(式30参照)
- (3)  $R_s R_s^T$  の固有ベクトルを計算し、そのうち最大の固有地に相当するベクトルを  $b_k$  とする。
- (4)  $R \leftarrow R - b_k b_k^T$  として  $R$  を更新する。



## アルゴリズム(2)

---

(5)  $k \leftarrow k+1$ とする。 $k > L$ ならばL個のベクトル $b_k$ が求められたことになるので処理を終了する。

そうでなければ(2)に戻る。

この結果、文書ベクトル $b_i$ を

$$\bar{d}_i = (b_1, \dots, b_L)^T d_i$$

のように次元の少ないL次元ベクトルに変換できる。  
このベクトルを使えばクラスタリングや特徴抽出を実行することが可能となる。