

# 第8章 次元縮約

## pLSIについて

7月1日  
鈴木直也

# 次元縮約

- 現実のデータ解析では、動もするとデータを表すベクトルが非常に高次元になります。
  - ベクトル間の距離が離れている。
  - 妥当な結果を得られない

次元が多くて困っているなら減らせばよい。

→ただ減らせばよいというものではなくできるだけ元の形を残して減らしていく。

# 射影

$n$ 次元空間の点を $k$ 次元空間に変換する処理を射影といいます。ここで $n$ 次元のデータが $m$ 個あるとき、このデータセットは行をデータに対応させる $m \times n$ の行列 $X$ で表現できます。

データを $k$ 次元に射影する1つの方法が、 $n \times k$ の行列 $A$ を掛けることです。このとき $XA$ の行ベクトルが $k$ 次元に縮約されたデータのベクトルになります。次元縮約を行うには多数ある $A$ の中から、何らかの条件を満たした $A$ を求める必要があります。

# 特異値分解

次元縮約を行う標準的な手法です。特異値分解では  $m \times n$  の行列  $X$  を3つの行列  $U, \Sigma, V^t$  の積に分解します。

The diagram illustrates the SVD equation:  $X = U \Sigma V^t$ . Matrix  $X$  is shown as a blue rectangle with dimensions  $m \times n$  below it. Matrix  $U$  is a green rectangle with dimensions  $m \times r$  below it. Matrix  $\Sigma$  is a smaller green square with dimensions  $r \times r$  below it. Matrix  $V^t$  is a green rectangle with dimensions  $r \times n$  below it. The matrices are connected by multiplication symbols ( $\times$ ) and an equals sign ( $=$ ).

図8.1: 特異値分解

$\Sigma$  は対角行列で、 $i$  行  $i$  列の対角要素を  $\lambda_i$  と置くと  $\lambda_1 \geq \lambda_2 \geq \lambda_3 \geq \dots \lambda_r \geq 0$  の関係があります。

ここで行列 $V$ を右から掛けると

$$\begin{array}{c} \boxed{X} \\ m \times n \end{array} \times \begin{array}{c} \boxed{V} \\ n \times r \end{array} = \begin{array}{c} \boxed{U} \\ m \times r \end{array} \times \begin{array}{c} \boxed{\Sigma} \\ r \times r \end{array}$$

となります。

ここで $XV$ は $m \times r$ の行列であり、元の $m \times n$ の行列が縮約されたこととなります。

# 次元縮約の結論

- $U$ の列ベクトルは $X$ の列ベクトルの張る空間の正規直行基底である。
- $V^t$ の行ベクトルは $X$ の行ベクトルの張る空間の正規直行基底である。
- $\lambda_i$ は $U$ の  $i$  番目の列ベクトル(あるいは $V^t$ の  $i$  番目の行ベクトル)の基底としての重要度を表す。

# LSI

これは次元縮約の手法の1つですが、特異値分解そのものです。ベクトル空間モデルで表現された文書ベクトルを特異値分解によって次元縮約する場合にLSIと呼びます。

文書ベクトルの各要素に重みをつける必要がある

→理論的な弱点

# pLSI

文書ベクトルを想定した次元縮約の手法。  
クラスタリングに直接利用することもできます。

文書ベクトルに重みも必要なく、頻度を要素とするベクトルでかまわない。

# Aspectモデル

文書と単語を結びつける潜在的なクラスを想定したモデル。  
文書 $d$ と単語 $w$ の出現をクラス $z$ を用いて以下のようにモデル化  
をする。

$$p(w|d) = \sum_z p(w|z)p(z|d)$$

ベイズの定理から

$$p(z|d) = \frac{p(z)p(d|z)}{p(d)}$$

また、 $p(d, w) = p(d)p(w|d)$  なので

# 続き

$$p(d, w) = \sum_z p(z) p(w|z) p(d|z) \quad (8.1)$$

となります。

K次に縮約する場合クラスをK個に設定します。

$$z_1, z_2, z_3, \dots, z_K$$

文書dに対して  $(p(z_1, d), p(z_2, d), \dots, p(z_K, d))$

が縮約されたベクトルとなります。

潜在的なクラスをそのままクラスタリングにおけるクラスと捉える。

→次元縮約の結果がクラスタリングを表す。

つまり、以下の式でデータdのクラスタ番号が得られます。

$$\arg \max_K p(d, z_k) = \arg \max_K p(z_k) p(d | z_k)$$

$p(z), p(w|z), p(d|z)$ を求めるにおいて一般的な式を求める必要はない。

与えられた文書集合が  $D = \{d_1, d_2, \dots, d_N\}$  でそれに関する単語の集合が  $W = \{w_1, w_2, \dots, w_m\}$  である場合、各  $k$  に対する  $p(z_k)$ 、各  $k$  と  $m$  に対する  $p(w_m | z_k)$ 、各  $k$  と  $n$  に対する  $p(d_n | z_k)$  が求まればよい。

→全部で  $K(1+M+N)$  個のパラメータを求める。

# 未知数(パラメータ)

パラメータは最尤法で求めることができる。

文書 $d$ に含まれる単語 $w$ の数を $n(d, w)$  で表すと対数尤度関数 $L$ は以下になります。

$$L = \sum_{i=1}^N \sum_{j=1}^M n(d_i, w_j) \log p(d_i, w_j)$$

式8.1より

$$L = \sum_{i=1}^N \sum_{j=1}^M n(d_i, w_j) \log p \left( \sum_{k=1}^K p(z_k) p(w_i | z_k) p(d_i | z_k) \right) \quad (8.2)$$

Lを最大化するパラメータを求める

→EMアルゴリズム<sup>\*1</sup>を用いる。

ここで尤度関数に隠れ変数 $z$ が  
埋め込まれているのでQ関数は  
直接求められる。

---

\*1: 混合分布モデル参照

# zの分布

z以外のパラメータを固定したときのzの分布を求める。

→つまり  $p(z | d, w)$

Aspectモデルでは  $p(d, w, z) = p(z)p(w | z)p(d | z)$

が定義されているので

$$p(z_k | d, w) = \frac{p(d, w, z_k)}{p(d, w)} = \frac{p(z_k)p(w | z_k)p(d | z_k)}{\sum_{k=1}^K p(z_k)p(w | z_k)p(d | z_k)}$$

# 続き

簡略化のために  $p(z_k | d_i, w_j) = Q_{ijk}$  と置くと

$$Q_{ijk}^{(t+1)} = \frac{p(z_k)^{(t)} p(w | z_k)^{(t)} p(d | z_k)^{(t)}}{\sum_{k=1}^K p(z_k)^{(t)} p(w | z_k)^{(t)} p(d | z_k)^{(t)}}$$

となります。

# zを固定した場合

zを固定つまり  $Q_{ijk}$  を固定したときのその他のパラメータを求める。

$$L = \sum_{i=1}^N \sum_{j=1}^M n(d_i, w_j) \log \left( \sum_{k=1}^K p(z_k) p(w_j | z_k) p(d_i | z_k) \right) \quad (8.3)$$

$$= \sum_{i=1}^N \sum_{j=1}^M n(d_i, w_j) \log \left( \sum_{k=1}^K Q_{ijk} \frac{p(z_k) p(w_j | z_k) p(d_i | z_k)}{Q_{ijk}} \right) \quad (8.4)$$

$$\leq \sum_{i=1}^N \sum_{j=1}^M n(d_i, w_j) \sum_{k=1}^K Q_{ijk} \log \frac{p(z_k) p(w_j | z_k) p(d_i | z_k)}{Q_{ijk}} \quad (8.5) * 1$$

---

\* 1 : Jensenの不等式より

# 続き

さらに変形して

$$L = \sum_{i=1}^N \sum_{j=1}^M n(d_i, w_j) \sum_{k=1}^K Q_{ijk} \log \left( p(z_k) p(w_j | z_k) p(d_i | z_k) \right) \\ - \sum_{i=1}^N \sum_{j=1}^M n(d_i, w_j) \sum_{k=1}^K Q_{ijk} \log Q_{ijk}$$

第2項は定数になる。

→ 第1項を再びLと置く

# ラグランジュの未定乗数法の利用

前式よりLの最大化は

$$\begin{cases} \sum_{i=1}^N p(d_i, z_k) = 1 \\ \sum_{j=1}^M p(w_j, z_k) = 1 \\ \sum_{k=1}^K p(z_k) = 1 \end{cases}$$

の関係があるのでラグランジュの未定乗数法を利用してとける。

$$L + \sum_{k=1}^K \alpha_k \left( 1 - \sum_{i=1}^N p(d_i | z_k) \right) + \sum_{k=1}^K \beta_k \left( 1 - \sum_{j=1}^M p(w_j | z_k) \right) + \gamma \left( 1 - \sum_{j=1}^M p(z_k) \right)$$

$$p(d_i | z_k) = u_{ik}, p(w_j | z_k) = v_{jk}, p(z_k) = w_k$$

とにおいて上記の式を  $u_{ik}, v_{jk}, w_k$   
で偏微分して極値問題を解く

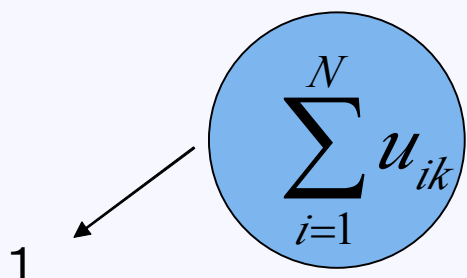
まず、 $u_{ik}$  を解きます。

$$\frac{\sum_{j=1}^M n(d_i, w_j) Q_{ijk}}{u_{ik}} - \alpha_k = 0$$

よって

$$u_{ik} = \frac{\sum_{j=1}^M n(d_i, w_j) Q_{ijk}}{\alpha_k}$$

両辺を*i*に関して和をとると

$$\sum_{i=1}^N u_{ik} = \frac{\sum_{i=1}^N \sum_{j=1}^M n(d_i, w_j) Q_{ijk}}{\alpha_k}$$


よって前式は

$$\alpha_k = \sum_{i=1}^N \sum_{j=1}^M n(d_i, w_j) Q_{ijk}$$

以上より

$$u_{ik} = p(d_i, z_k) = \frac{\sum_{j=1}^M n(d_i, w_j) Q_{ijk}}{\sum_{i=1}^N \sum_{j=1}^M n(d_i, w_j) Q_{ijk}}$$

更新式の形で書くと

$$p(d_i, z_k)^{(t)} = \frac{\sum_{j=1}^M n(d_i, w_j) Q_{ijk}^{(t)}}{\sum_{i=1}^N \sum_{j=1}^M n(d_i, w_j) Q_{ijk}^{(t)}}$$

同様に  $\nu_{jk}, w_k$  でそれぞれ偏微分を行うと以下の式を得られます。

$$p(w_j, z_k)^{(t)} = \frac{\sum_{i=1}^N n(d_i, w_j) Q_{ijk}^{(t)}}{\sum_{i=1}^N \sum_{j=1}^M n(d_i, w_j) Q_{ijk}^{(t)}}$$

$$p(z_k)^{(t)} = \frac{\sum_{i=1}^N \sum_{j=1}^M n(d_i, w_j) Q_{ijk}^{(t)}}{\sum_{k=1}^K \sum_{i=1}^N \sum_{j=1}^M n(d_i, w_j) Q_{ijk}^{(t)}}$$