

機械学習ゼミ

⑤ 共クラスタリング

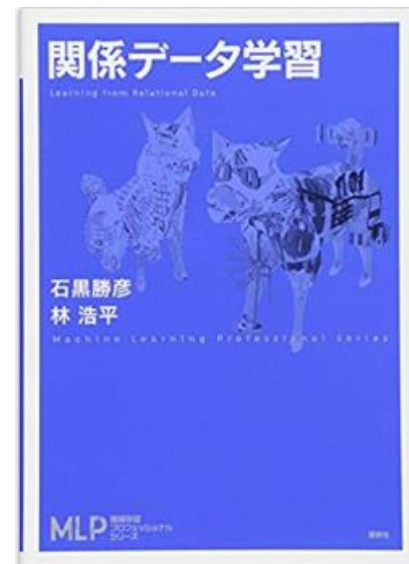
2018年7月7日(土)

五百藏夏穂

関係データ学習

- * スペクトラルクラスタリング
- * 確率的ブロックモデル/無限関係モデル
- * 行列分解

本ゼミでの説明は、
“関係データ学習” 石黒勝彦, 林浩平(著)
に準拠します。



関係データとは

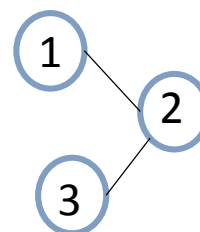
関係データ 複数のデータの中に観測, 定義される関係に着目したデータのこと (SNSユーザー間, 購買データ, etc)

有向関係
無向関係

起点によって関係が変わるか否か

単一ドメイン
複数ドメイン

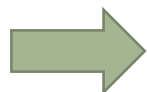
オブジェクトが同じ種類か否か
ex (単一: SNSユーザー関係, 複数: 購買データ)



グラフ表現
 $V = \{1, 2, 3\}$

0	1	0
1	0	1
0	1	0

行列表現
 $X = (x_{i,j}) \ i, j = 1, 2, 3$



対称関係データ: 単一ドメインかつ無向関係
スペクトラルクラスタリング

非対称関係データ: それ以外

確率的ブロックモデル/無限関係モデル, 行列分解

タスク

- * 予測 : 関係データ学習を通してみ観測のデータの値を推定
→ アイテム推薦問題
- * 知識抽出 : 関係データの特徴を解析, モデル化し, 知見や知識の情報抽出
→ コミュニティ抽出, クラスタリング

タスク

関係データの種類

		知識抽出	予測
二項関係	対称	スペクトラル クラスタリング	行列分解
	非対称	確率的ブロックモデル 無限関係モデル	
多項関係			テンソル分解

関係データ学習におけるクラスタリング

対象: オブジェクト(ノード) V のクラスタリング

類似性: 各オブジェクト i, j 間の“関係”に基づく分類

※対象: エッジ(リンク), 類似性: オブジェクト自身が持つ特徴量に基づくクラスタリングはここでは本ゼミでは対象外とします。

教師あり	・分類 (決定木, NN, ...etc)
	・回帰分析 (線形回帰, ...etc)
機械学習	・クラスタリング (k-means, ...etc)
	・共起分析 (協調フィルタリング, ...etc) ・次元圧縮 (主成分分析, 特異値分解, ...etc)
教師なし	・Q学習
	・TD学習
強化学習	
etc	

コミュニティ検出と密結合グラフ

密結合クラスタ: クラスタ内部で密な結合, クラスタ外部と疎な結合

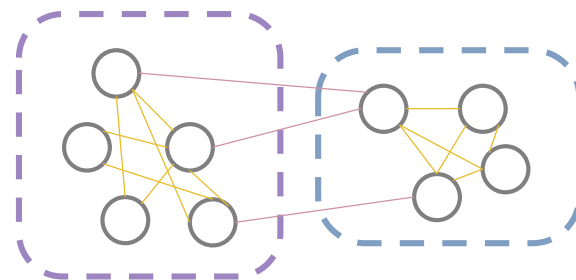
疎結合クラスタ: クラスタ内部で疎な結合, クラスタ外部と疎な結合



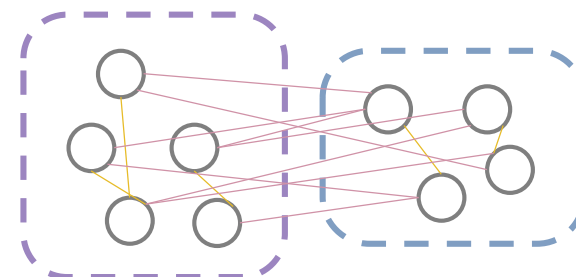
コミュニティ検出では前者

どのクラスタにどのオブジェクトを帰属させるかを計算

密結合クラスタグラフ



疎結合クラスタグラフ



グラフカット

疎・密結合の発見や判断に用いる特徴量 **グラフカット**

$$\text{cut}(P) = \frac{1}{2} \sum_{k=1}^K \sum_{i \in P_k, j \in \bar{P}_k} x_{i,j}$$

P_k : k 番目クラスタに所属する頂点集合

\bar{P}_k : k 番目クラスタに所属しない頂点集合

$x_{i,j}$: 異なるクラスタに所属する頂点 i, j 間の重み

頂点間のつながりの弱い部分で分割すると考えて

➡ **カット最小化問題**として密結合クラスタを抽出

スペクトラルクラスタリング

グラフラプラシアン固有値に着目し、
固有値の値によってグラフのノード間の連結性が推定できる



固有値分解とk-meansの併用でグラフを構成する
頂点のクラスタリングを行う手法

①非正規化グラフラプラシアンによるスペクトラルクラスタリング

グラフラプラシアン

観測データ行列:隣接行列

次数行列: $D = (d_{i,j}), \quad i, j = 1, 2, \dots, N$

$$d_{i,j} = \begin{cases} \sum_{j=1}^N x_{i,j} & i = j \\ 0 & i \neq j \end{cases}$$

グラフラプラシアン:

$$L = X - D$$

観測データ行列X

	j=1	j=2	j=3	j=4	j=5
i=1	0	1	0	1	0
i=2	1	0	0	1	0
i=3	0	0	0	0	1
i=4	1	1	0	0	0
i=5	0	0	1	0	0

次数行列D

	j=1	j=2	j=3	j=4	j=5
i=1	2	0	0	0	0
i=2	0	2	0	0	0
i=3	0	0	1	0	0
i=4	0	0	0	2	0
i=5	0	0	0	0	1

グラフラプラシアンL

	j=1	j=2	j=3	j=4	j=5
i=1	2	-1	0	-1	0
i=2	-1	2	0	-1	0
i=3	0	0	1	0	-1
i=4	-1	-1	0	2	0
i=5	0	0	-1	0	1

入力: 対称データ行列X, クラスタ数K

次数行列Dの計算

グラフラプラシアンLの計算

Lの固有値分解を計算
固有値を昇順に並び替えて最初K個の
列固有ベクトルを $V = (v_1, \dots, v_K) \in R^{N \times K}$
として取り出す

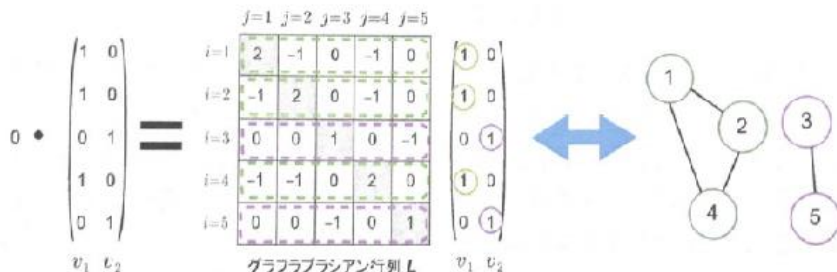
行列Vの各行ベクトルを
k-meansでクラスタリング

出力: K次元空間オブジェクト配置行列V
オブジェクトのクラスタ割り当てZ

①非正規化グラフラプラシアンによるスペクトラルクラスタリング

固有値分解によるクラスタ抽出

$$L = V\Lambda V^T \quad \text{または} \quad \lambda_k v_k = L v_k$$



L は半正定値行列なので

$$0 \leq \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_K$$

この時定数ベクトルの固有値は0なので最小値 $\lambda_1=0$ を取る

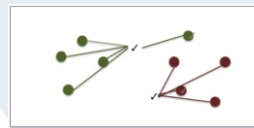
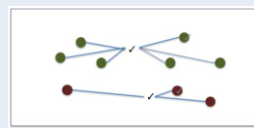
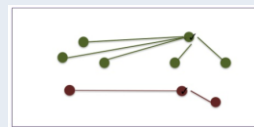
定数ベクトルの要素を分割したベクトルを取ることで、 K 個のクラスタを持つグラフでは最小 K 個の固有値0 (理想状態)

クラスタ割り当ての計算

列固有ベクトル $V = (v_1, \dots, v_K)$ を N 個の K 次元ベクトルととらえてk-meansを適用

クラスタ割り当て $Z = \{z_i = k\}$ が求まる

※補足 : k-means



K 個の中心を決める

↓

最も近くの中心クラスタに割り当て

↓

クラスタ内平均からクラスタ中心を更新

↓

中心が変化しなくなったら終了

②正規化グラフラシアンによるスペクトラルクラスタリング

正規化カット

カット最小化の傾向として
大きなクラスタと、孤立し多数の単頂点クラスタ
ができやすい



各クラスタの大きさに制約をかける
ように正規化を行ったカットを用いる

レシオカット

$$\text{RatioCut} = \frac{1}{2} \sum_{k=1}^K \frac{\sum_{i \in P_k, j \in \bar{P}_k} x_{i,j}}{|P_k|} = \sum_{k=1}^K \frac{\text{cut}(P_k, \bar{P}_k)}{|P_k|}$$

- ・各クラスタのサイズ(頂点数)で割り正規化
→クラスタが大きいとカットが小さくなる
- ・非正規化グラフラシアンのアルゴリズムで解ける

ノーマライズドカット

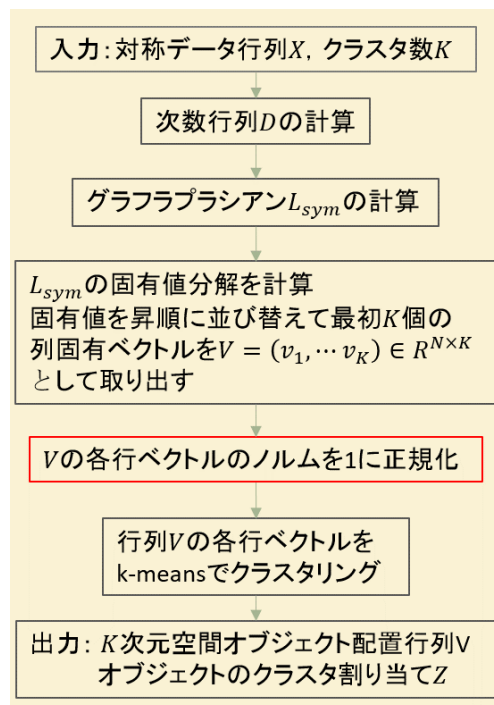
$$\text{NCut} = \frac{1}{2} \sum_{k=1}^K \frac{\sum_{i \in P_k, j \in \bar{P}_k} x_{i,j}}{\text{vol}(P_k)} = \sum_{k=1}^K \frac{\text{cut}(P_k, \bar{P}_k)}{\text{vol}(P_k)} \quad \text{vol}(P_k) = \sum_{i,j \in P} x_{i,j}$$

- ・クラスタ内のリンクの重みの総和で割る
→クラスタが大きかつ相互に密に結合しているとカットが小さくなる
- ・正規化グラフラシアンに基づくスペクトラルクラスタリングに帰着

②正規化グラフラプラシアンによるスペクトラルクラスタリング

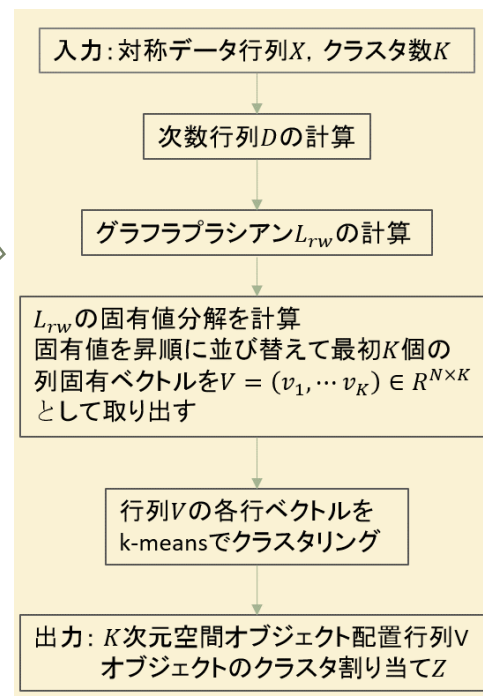
対称正規化ラプラシアン

$$L_{sym} = D^{-\frac{1}{2}} L D^{-\frac{1}{2}}$$



酔歩正規化ラプラシアン

$$L_{rw} = D^{-1} L$$



比較

* コミュニティ(密結合)が対象であれば、密結合も考慮したノーマライズドカットの方が望ましい

* L_{sym} はk-means前の正規化のため、計算コストが大きい。そのため、 L_{rw} の方が望ましい

ノーマライズドカットをグラフラプラシアンを用いた二次形式に書き換え、それを最小化する解を制約付きで解くと導出できる

関係データXの重みに従って酔歩すると想定すると遷移確率行列は $T = D^{-1} X$
ノーマライズドカット最小化 = クラスタ間遷移確率最小化

②実データ適用例

Zachary's Karate Club network data

空手クラブのメンバーの交友関係は、レッスン料の違いから2つの集団に分裂

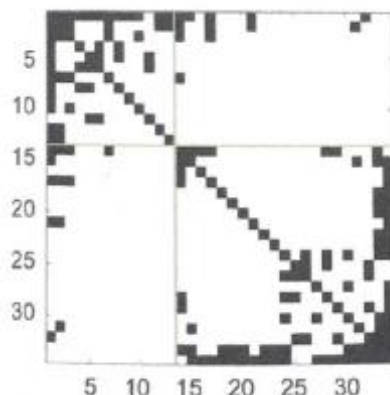
$N=34$, $K=2$ として各手法を適用

実際の結果と一致しなかったもの

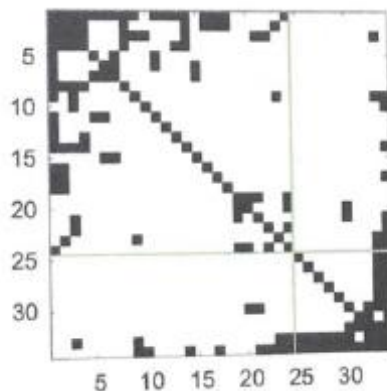
非正規化/酔歩正規化 : 4/34

対称正規化 : 7/34

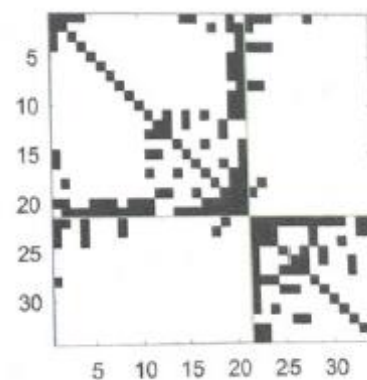
非正規化



対称正規化



酔歩正規化



確率的ブロックモデルと無限関係モデル

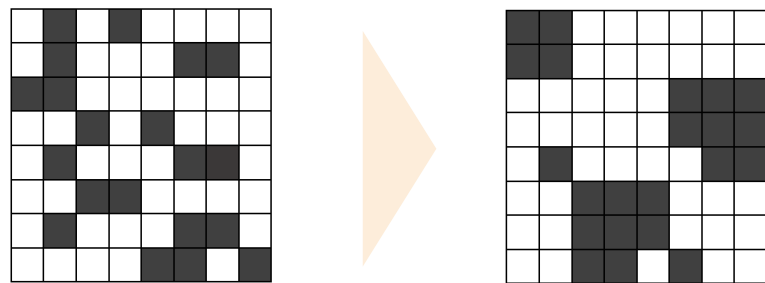
スペクトラルクラスタリングでは、対象関係データの密結合クラスタに対象が限られていた

➡ 非対称関係データにも適応可能で
密結合以外も抽出可能なデータクラスタリング手法

* 潜在的なブロック構造を仮定

* 確率的アプローチ

関係データに潜むブロック構造を推定し、この時、
この構造によりどのような観測データが確率的に得られるのか？



確率的ブロックモデル
Stochastic Block Model (SBM)

無限関係モデル
Infinite Relational Model (IRM)

クラスタ数Kが
自動的に求まる

確率的生成モデル SBMもIRMもこの一種である

系を構成する各変数の値が依存関係に従って順番に生成される過程を確率で表現するモデル

ex)サイコロの目(100回振って出た目を記録)

100回サイコロを振る : 変数をサンプリング
1回ごとにサイコロの形状に従って : パラメータで規定される確率分布
ランダムに目の値が決まる : 観測値が生成される

$$x_i | \pi \sim P(\pi), i = 1, 2, \dots, 100$$

$X = (x_1, \dots, x_{100})$: 観測値
 $\pi = (\pi_1, \dots, \pi_d, \pi_D)$: 面の出やすさパラメータ
 D : 面数 $\sum_d \pi_d = 1$

- メリット
- * 確率分布の形状やパラメータが全て求めれば、実際に値を生成できる
 - * 事前知識を確率分布の形状や種類として反映できる

確率モデルのベイズ推定

観測データから確率分布の形状やパラメータ値を推定
本書で説明されるベイズ推定では、ベイズ事後分布を求める

$$\underline{p(Z, \{\theta\} | X)} \propto p(X | Z, \{\theta\}) p(Z, \{\theta\})$$

周辺化

$$\underline{p(Z | X)} \propto \int p(X | Z, \{\theta\}) p(Z, \{\theta\}) d\{\theta\}$$

確率的ブロックモデルの定式化

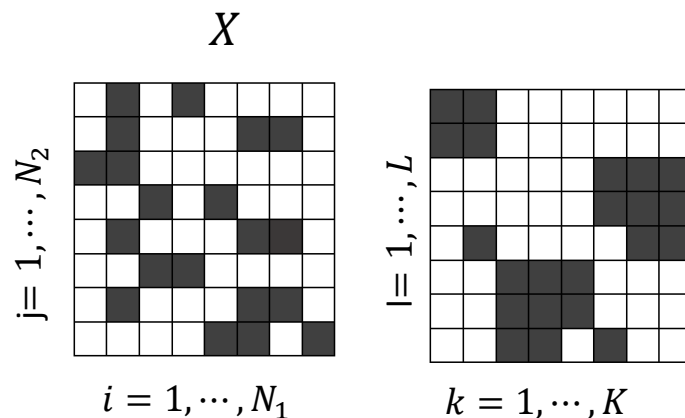
観測データ: $X = \{x_{i,j}\}$ $x_{i,j} \in \{1,0\}$

行(第一ドメイン)インデックス: $i = 1, \dots, N_1$

列(第二ドメイン)インデックス: $j = 1, \dots, N_2$

行クラスインデックス: $k = 1, \dots, K$

列クラスインデックス: $l = 1, \dots, L$



$$\pi_1 | \alpha_1 \sim \text{Dirichlet}(\alpha_1)$$

$$\pi_2 | \alpha_2 \sim \text{Dirichlet}(\alpha_2)$$

$$z_{1,i} = k | \pi_1 \sim \text{Discrete}(\pi_1)$$

$$z_{2,j} = l | \pi_2 \sim \text{Discrete}(\pi_2)$$

$$\theta_{k,l} | a_0, b_0 \sim \text{Beta}(a_0, b_0)$$

$$x_{i,j} | \{\theta_{k,l}\}, z_{1,i}, z_{2,j} \sim \text{Bernoulli}(\theta_{z_{1,i}, z_{2,j}})$$

π : オブジェクトのクラスタへの混合割合パラメータ

α : ハイパラメータ

ex) K面のサイコロをランダムに生成する分布

各ドメインのクラスタの混合割合に従ってオブジェクトのクラスタ割り当てをサンプリング

ex) K面のサイコロに従って目を出す

θ : クラスタ間関係強さ(2値関係なのでベータ分布)

z : 割り当てを表す隠れ変数

観測行列の生成

周辺化ギブスサンプラー(CGS)による推論

$$z_i^{(s)} \sim p(z_i | X, Z^{/(i)}) \propto \int p(X^{(i)} | X^{/(i)}, z_i, Z^{/(i)} \{\theta\}) p(z_i, \{\theta\} | X^{/(i)}, Z^{/(i)}) d\{\theta\}$$

※確率モデルは隠れ変数 z , パラメータ集合 $\{\theta\}$, 観測データ X からなるが、結果的に必要なのは隠れ変数の事後分布だけなのでパラメータは周辺化して推論

s : サンプル回数
 $X^{(i)}$: z_i に依存する観測量
 $X^{/(i)}$: $X - X^{(i)}$

事後分布から生成されたとみなされる確率変数のサンプリングを多数行い、サンプルの経験分布をもって事後分布の推定値を計算する
CGSでは、1回のサンプリングである一つの変数だけのサンプリングを順番に行う。

$z_i^{(s)}$ は s 回目の事後分布サンプリング結果



$z_i \leftarrow z_i^{(s)}$: 既知の値として保存

z_j ($i \neq j$) のサンプリングに映る

CGSによるSBMの推論式

第一・第二ドメインで対象なので

$$p(z_{1,i} = k | X, Z_1^{(i)}, Z_2) \propto$$

$$\frac{\int p\left(X^{(i)} | X^{/(i)}, z_{1i=k}, Z_1^{(i)}, Z_2, \pi_1, \pi_2, \{\theta_{k,l}\}\right) \cdot p\left(z_{1i=k}, \pi_1, \pi_2, \{\theta_{k,l}\} | X^{/(i)}, Z_1^{(i)}, Z_2\right) d\pi_1 d\pi_2 d\{\theta_{k,l}\}}{z_{1i} \text{に依存する関係データ行列の要素に対する尤度} \quad z_{1i} \text{及びパラメータの} z_{1i} \text{以外の情報が所与の時の事後分布}}$$

$$p(z_{1,i} = k | X, Z_1^{(i)}, Z_2) \propto$$

$$\frac{\int p(z_{1,i} = k | \pi_1) p(\pi_1 | Z_1^{\bar{i}}) d\pi_1 \cdot \int \prod_{l=1}^L \prod_{j=1}^{N_2} [p(x_{i,j} | \theta_{k,l})]^{\Pi(z_{2,j}=l)} p(\theta_{k,l} | X^{/(i)}, Z_1^{(i)}, Z_2) d\{\theta_{k,l}\}}{\pi_1 \text{を周辺化した時のクラスタ割り当て} z_{1,i} \text{の確率} \quad \text{第二ドメインのクラスタごと} \text{にクラスタ間関係強さ} \theta_{k,l} \text{を} \text{周辺化した時の観測データの尤度の影響}}$$

生成モデルに従って計算

$$p(z_{1,i} = k | X, Z_1^{(i)}, Z_2) \propto \hat{\alpha}_{1,k} \frac{\Gamma(\hat{a}_{k,l} + \hat{b}_{k,l})}{\Gamma(\hat{a}_{k,l})\Gamma(\hat{b}_{k,l})} \times \frac{\Gamma(\hat{a}_{k,l} + \sum_{j=1}^{N_2} x_{i,j} \Pi(z_{2,j}=l)) \Gamma(\hat{b}_{k,l} + \sum_{j=1}^{N_2} (1-x_{i,j}) \Pi(z_{2,j}=l))}{\Gamma(\hat{a}_{k,l} + \hat{b}_{k,l} + \sum_{j=1}^{N_2} \Pi(z_{2,j}=l))}$$

$$\hat{a}_{k,l} = a_0 + \hat{n}_{k,l}^{(+)} \quad \hat{b}_{k,l} = b_0 + \hat{n}_{k,l}^{(-)}$$

CGSによるSBMの推論式(十分統計量)

定義

$$m_{1,k} = \sum_{i=1}^{N_1} \Pi(z_{1,i} = k) \quad \text{: クラスタ } k \text{ に所属するオブジェクト数}$$

$$n_{k,l}^{(+)} = \sum_{i=1}^{N_1} \sum_{j=1}^{N_2} x_{i,j} \Pi(z_{1,i} = k) \Pi(z_{2,j} = l) \quad \text{: ブロック } (k, l) \text{ において } x = 1(0) \text{ となる要素数}$$

$$n_{k,l}^{(-)} = \sum_{i=1}^{N_1} \sum_{j=1}^{N_2} (1 - x_{i,j}) \Pi(z_{1,i} = k) \Pi(z_{2,j} = l)$$

$z_{1,i}$ のサンプリング時に

$z_{1,i}$ のクラスタ割り当てを解除してた時の一つ抜き十分統計量

$$\hat{m}_{1,k} = \sum_{i' \neq i, i'=1}^{N_1} \Pi(z_{1,i'} = k) = m_{1,k} - \Pi(z_{1,i} = k)$$

$$\hat{n}_{k,l}^{(+)} = \sum_{i' \neq i, i'=1}^{N_1} \sum_{j=1}^{N_2} x_{i',j} \Pi(z_{1,i'} = k) \Pi(z_{2,j} = l) = n_{k,l}^{(+)} - \Pi(z_{1,i} = k) \sum_{j=1}^{N_2} x_{i,j} \Pi(z_{2,j} = l)$$

$$\hat{n}_{k,l}^{(-)} = \sum_{i' \neq i, i'=1}^{N_1} \sum_{j=1}^{N_2} (1 - x_{i',j}) \Pi(z_{1,i'} = k) \Pi(z_{2,j} = l) = n_{k,l}^{(-)} - \Pi(z_{1,i} = k) \sum_{j=1}^{N_2} (1 - x_{i,j}) \Pi(z_{2,j} = l)$$

更新

$$m_{1,k} = \hat{m}_{1,k} + \Pi(z_{1,i} = k)$$

$$n_{k,l}^{(+)} = \hat{n}_{k,l}^{(+)} + \Pi(z_{1,i} = k) \sum_{j=1}^{N_2} x_{i,j} \Pi(z_{2,j} = l)$$

$$n_{k,l}^{(-)} = \hat{n}_{k,l}^{(-)} + \Pi(z_{1,i} = k) \sum_{j=1}^{N_2} (1 - x_{i,j}) \Pi(z_{2,j} = l)$$

対象となる変数の差し引きだけで、1回のサンプリングに対応する更新式計算ができる！

無限関係モデルの定式化

無限個のクラスタの存在を仮定した確率的生成モデル
潜在クラスタ数 K, L を自動的に決定できる

$$Z_1 | \alpha_1 \sim \text{CRP}(\alpha_1)$$

$$Z_2 | \alpha_2 \sim \text{CRP}(\alpha_2)$$

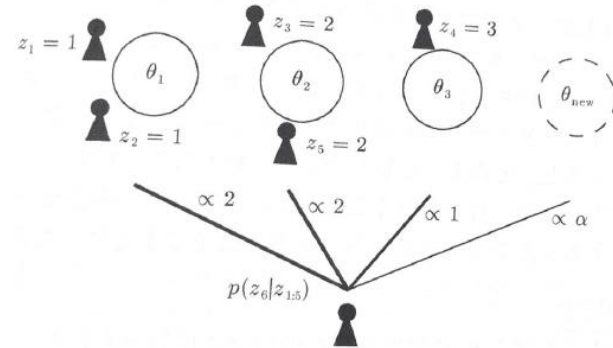
CRPは無次元ディリクレ分布と無次元離散分布を組み合わせたもの

$$\theta_{k,l} | a_0, b_0 \sim \text{Beta}(a_0, b_0)$$

$$x_{i,j} | \{\theta_{k,l}\}, z_{1,i}, z_{2,j} \sim \text{Bernoulli}(\theta_{z_{1,i}, z_{2,j}})$$

Chinese Restaurant Process (CRP)

ノンパラメトリックベイズで用いられる確率過程



- * 各テーブルに着席できる客数に上限はない
- * 中華料理店に配置できるテーブル数に上限はない

➡ 入店した客をどのテーブルに配置するか決める

$n + 1$ 人目の客のテーブル選択確率

$$P(z_{n+1} = k | z_{1:n}, \alpha) \propto \begin{cases} \sum_{i=1}^n \mathbb{I}(z_i = k) & k \in \{1, 2, \dots, K\} \\ \alpha & k = K + 1 \end{cases}$$

無限個のコンポーネントを持つ混合分布を正しく表現
実際のサンプリングでは $N+1$ 個のテーブルと、有限になるので実装可能

CGSによるIRMの推論式

$$p(z_{1,i} = k | X, Z_1^{/(i)}, Z_2) \propto \int p\left(X^{(i)} | X^{/(i)}, z_{1i=k}, Z_1^{/(i)}, Z_2, \{\theta_{k,l}\}\right) \cdot p\left(z_{1i=k}, \{\theta_{k,l}\} | X^{/(i)}, Z_1^{/(i)}, Z_2\right) d\{\theta_{k,l}\}$$

$$p(z_{1,i} = k | X, Z_1^{/(i)}, Z_2) \propto p\left(z_{1,i} = k | Z_1^{/(i)}\right) \int \prod_{l=1}^L \prod_{j=1}^{N_2} [p(x_{i,j} | \theta_{k,l})]^{\Pi(z_{2,j}=l)} p(\theta_{k,l} | X^{/(i)}, Z_1^{/(i)}, Z_2) d\{\theta_{k,l}\}$$

場合分けして計算

$$p(z_{1,i} = k \in \{1, 2, \dots, K\} (\text{既存クラスター}) | X, Z_1^{/(i)}, Z_2) \propto \hat{m}_{1,k} \prod_{l=1}^L \left[\frac{\Gamma(\hat{a}_{k,l} + \hat{b}_{k,l})}{\Gamma(\hat{a}_{k,l})\Gamma(\hat{b}_{k,l})} \times \frac{\Gamma(\hat{a}_{k,l} + \sum_{j=1}^{N_2} x_{i,j} \Pi(z_{2,j}=l)) \Gamma(\hat{b}_{k,l} + \sum_{j=1}^{N_2} (1-x_{i,j}) \Pi(z_{2,j}=l))}{\Gamma(\hat{a}_{k,l} + \hat{b}_{k,l} + \sum_{j=1}^{N_2} \Pi(z_{2,j}=l))} \right]$$

$$p(z_{1,i} = k + 1 (\text{新規クラスター}) | X, Z_1^{/(i)}, Z_2) \propto \alpha_1 \prod_{l=1}^L \left[\frac{\Gamma(a_0 + b_0)}{\Gamma(a_0)\Gamma(b_0)} \times \frac{\Gamma(a_0 + \sum_{j=1}^{N_2} x_{i,j} \Pi(z_{2,j}=l)) \Gamma(b_0 + \sum_{j=1}^{N_2} (1-x_{i,j}) \Pi(z_{2,j}=l))}{\Gamma(a_0 + b_0 + \sum_{j=1}^{N_2} \Pi(z_{2,j}=l))} \right]$$

CGSによるIRMの推論式(十分統計量)

更新

$$m_{1,k} = \hat{m}_{1,k} + \Pi(z_{1,i} = k)$$
$$n_{k,l}^{(+)} = \hat{n}_{k,l}^{(+)} + \Pi(z_{1,i} = k) \sum_{j=1}^{N_2} x_{i,j} \Pi(z_{2,j} = l)$$
$$n_{k,l}^{(-)} = \hat{n}_{k,l}^{(-)} + \Pi(z_{1,i} = k) \sum_{j=1}^{N_2} (1 - x_{i,j}) \Pi(z_{2,j} = l)$$

もし $z_{1,i} = k + 1$ ならば
ハットつきの値を全て0にする

出力方法

サンプリング実現値を用いて事後分布を計算

最初の適当なB期間のサンプルを棄却

その後M個置き of サンプル実現値だけを抽出 ($T = \frac{S-B}{M}$ 個)

$$p(z_{1,i} = k|X) = \frac{1}{T} \sum_{t=1}^T \Pi(z_{1,i}^{(t)} = k) \quad p(z_{2,j} = l|X) = \frac{1}{T} \sum_{t=1}^T \Pi(z_{2,j}^{(t)} = l)$$

具体的なクラス割り当て

事後分布最大化(MAP)に基づく

$$c_{1,i} = \operatorname{argmax} p(z_{1,i} = k|X)$$

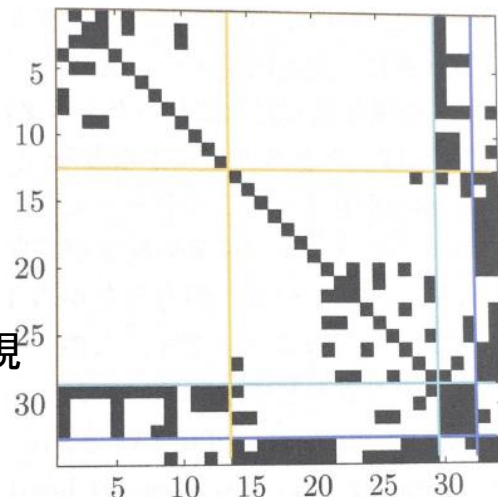
$$c_{2,j} = \operatorname{argmax} p(z_{2,j} = l|X)$$

②実データ適用例

* Zachary's Karate Club network data

4つに分裂

現実と分裂結果の一致度は
スペクトラルクラスタリングより良い値に.



人間の思い込みに反するがよりよく関係データを表現

* Enron社内のEメールデータ

メール送信者(K), 受信者(L)の二値データ

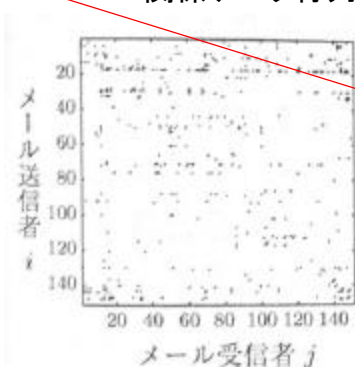
K=6, L=3クラスタに分割

経営層コミュニティ

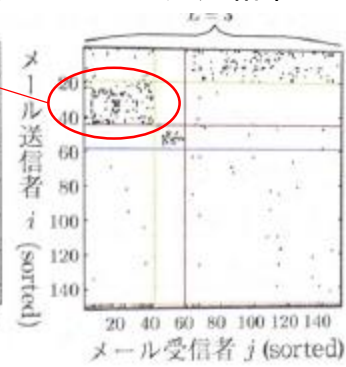


真のクラスタ構造がわからないデータでもクラスタリングを実行可能

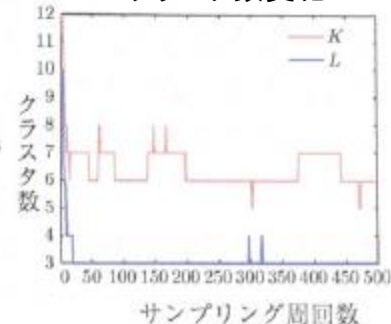
関係データ行列



IRM適用結果



繰り返し回数ごとのクラスタ数変化



単純行列分解

例として映画の推薦を考えます

I : 顧客 J : 映画

要素の値が評価値となる行列

$$I \begin{matrix} & J \\ & X \end{matrix} \cong I \begin{matrix} R \\ U \end{matrix} \begin{matrix} J \\ V^T \\ R \end{matrix}$$

$$x_{ij} \cong u_{i,action} v_{j,action} + u_{i,horror} v_{j,horror}$$

より一般に $R \in N$ 個の評価軸があり、すべての顧客、映画について評価値が得られているとする

X のランク R 行列分解

$$X \cong U V^T$$

因子行列

目的関数

$$E = X - U V^T$$

X と $U V^T$ の近似誤差を最小化する問題ととらえる

2乗誤差最小化問題



$$\min_{U,V} \frac{1}{2} \|E\|_{\text{Fro}}^2 = \min_{U,V} \frac{1}{2} \|X - U V^T\|_{\text{Fro}}^2$$

$u_{i,:} \in \mathbb{R}$: 顧客 i の嗜好をまとめたもの

$v_{j,:} \in \mathbb{R}$: 映画 j の成分をまとめたもの

$U \in \mathbb{R}^{I \times R}$: 行方向にまとめた

$V \in \mathbb{R}^{J \times R}$: 列方向にまとめた

l^2 正則行列分解

X に大きいノイズが加わっていた時などは過学習の危険がある
→正則化(制約を付加し, 解の範囲を狭める)

$$\min_{U,V} \frac{1}{2} \|X - UV^T\|_{\text{Fro}}^2 + \frac{\lambda}{2} (\|U\|_{\text{Fro}}^2 + \|V\|_{\text{Fro}}^2)$$

$\lambda \geq 0$ は正則化強さを決定する係数

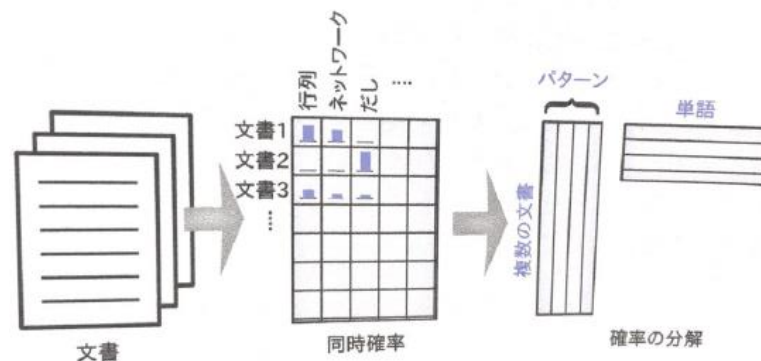
非負行列分解

テキストの単語数データなど, 非負値しかとらないもの

$$\min_{U \in \mathbb{R}_+^{I \times R}, V \in \mathbb{R}_+^{J \times R}} \frac{1}{2} \|X - UV^T\|_{\text{Fro}}^2$$

Ex) 同時確率
 $x_{i,j} = p(\text{文書}i, \text{単語}j)$

- * 単語の確率分布は典型的なパターンを持つ
- * 文書における単語の確率分布はパターンの組み合わせ



アルゴリズム

目的関数の導関数を0としても、連立方程式を代数的に解くことができない

➡ 勾配法による最適化

一次交互勾配法

$K \in N$ 個の変数 $\theta_1, \dots, \theta_K$ を引数にとる関数 $h(\theta_1, \dots, \theta_K) \in R$ では

$$\theta_k^{new} = \theta_k - \eta \nabla_{\theta_k} h(\theta_1^{new}, \dots, \theta_{k-1}^{new}, \theta_k, \theta_{k+1}, \dots, \theta_K) \quad \eta: \text{学習率}$$

l^2 正則行列分解の交互最適化

目的関数と連鎖率より

$$\nabla_U f_{l^2} = -XV + U(V^T V + \lambda I)$$

$$\nabla_V f_{l^2} = -X^T U + V(U^T U + \lambda I)$$

X が密でデータが大きい場合計算量の負荷が大きい

入力: 行列 X , ランク R , 正則化係数 λ ,
学習率 η , 精度 ϵ

U, V を乱数で初期化

繰り返し

$g \leftarrow f_{l^2}(U, V)$
 $U \leftarrow U - \nabla_U f_{l^2}(U, V)$
 $V \leftarrow V - \nabla_V f_{l^2}(U, V)$

until $|g - f_{l^2}(U, V)| / f_{l^2}(U, V) < \epsilon$

出力: 因子行列 U, V

疑似二次交互勾配法

一次交互勾配法は実装は比較的容易であるが、収束が遅いことも、
二階微分の情報を用いて収束を早くする

$$\theta_k^{new} = \theta - \eta (\nabla \nabla^T g(\theta))^{-1} \nabla g(\theta)$$

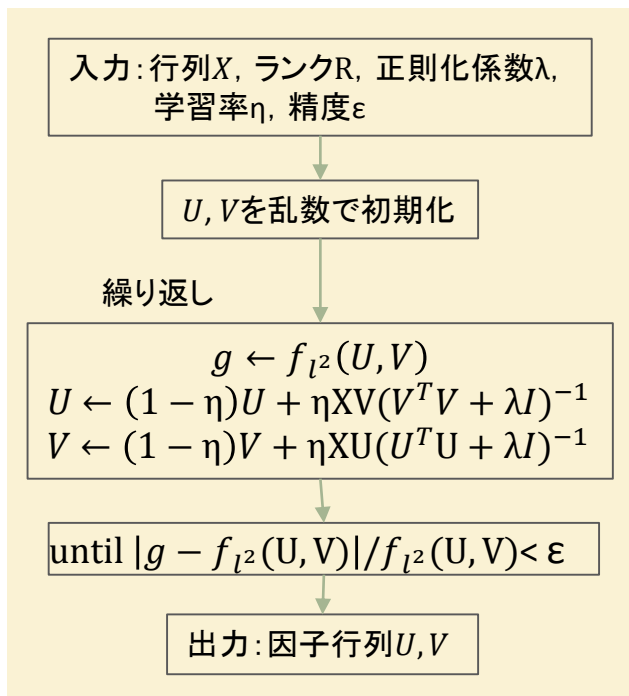
ヘッセ行列

一次勾配法に比べて収束速度は速いが、
ヘッセ行列の逆行列を計算しなければならず、
計算量は高次元

l^2 正則行列分解の交互最適化

$$U^{new} = (1 - \eta)U + \eta XV(V^T V + \lambda I)^{-1}$$

$$V^{new} = (1 - \eta)V + \eta XU(U^T U + \lambda I)^{-1}$$



交互勾配法＋射影勾配降下法(制約条件)

射影勾配降下法 非負制約などの制約を取り入れたい場合

u : 制約を満たす集合

$$U^{new} = \text{Proj}_u [U - \eta \nabla_U f_{l^2}(U, V)]$$
$$\text{Proj}_u [U] = \text{argmin}_{Z \in u} \|U - Z\|_{\text{Fro}}^2$$

通常の勾配法の後、空間 u からはみ出た部分に戻す

非負条件では、

$$U^{new} = [U - \eta \nabla_U f_{l^2}(U, V)]_+$$

$$[U]_+ : u_{i,j} > 0 \text{ の時 } u_{i,j}, u_{i,j} \leq 0 \text{ の時 } 0$$

入力: 行列 X , ランク R , 正則化係数 λ ,
学習率 η , 精度 ϵ

U, V を乱数で初期化

繰り返し

$$g \leftarrow f_{l^2}(U, V)$$
$$U \leftarrow [(1 - \eta)U + \eta XV(V^T V + \lambda I)^{-1}]_+$$
$$V \leftarrow [(1 - \eta)V + \eta XU(U^T U + \lambda I)^{-1}]_+$$

until $|g - f_{l^2}(U, V)| / f_{l^2}(U, V) < \epsilon$

出力: 因子行列 U, V

確率勾配降下法

X が密で I, J がお真理に大きいとメモリに乗りきらないことがある
少数のサンプルで勾配を近似的に計算し最適化を行う方法.

→ 目的関数 f がサンプル n のみに依存する関数 f_n
の和で書けるとき、代わりに f_n の勾配を使う

要素ごと

$$\nabla_{u_i} f_{l^2}^{(i,j)} = -x_{i,j} v_j + u_i (v_j^T v_j + \frac{\lambda}{J})$$

$$\nabla_{v_j} f_{l^2}^{(i,j)} = -x_{i,j} u_i + v_j (u_i^T u_i + \frac{\lambda}{I})$$

行列ごと

共起行列のような X が複数の行列の和で与えられており
要素へのアクセスにもコストがかかる場合 (ex. 単語の共起行列)

単純行列分解での目的関数を

$$f_{std}(X; U, V) = \frac{1}{2} \|X - UV^T\|_{Fro}^2 \text{と表記すると}$$

$$f_{std}(X; U, V) = P_m [f_{std}(MC^{(m)}; U, V)] + M^2 \sigma^2$$

$$\sigma^2 = P_m \|C^{(m)} - P_m[C^{(m)}]\|_{Fro}^2, P_m[\cdot] = \frac{1}{M} \sum_{m=1}^M (\cdot)$$

入力: 行列 X , ランク R , 正則化係数 λ ,
学習率 η , 反復回数 T

U, V を乱数で初期化

for $t = 1, \dots, T$

$[I]$ 上の一様乱数から i をサンプリング
 $[J]$ 上の一様乱数から j をサンプリング

$$u_i \leftarrow u_i + \eta(t) \nabla_{u_i} f_{l^2}^{(i,j)}(u_i, v_j)$$

$$v_j \leftarrow v_j + \eta(t) \nabla_{v_j} f_{l^2}^{(i,j)}(u_i, v_j)$$

end for

出力: 因子行列 U, V

入力: 行列 $\{C_1, \dots, C_M\}$, ランク R , 正則化係数 λ ,
学習率 η , 反復回数 T

U, V を乱数で初期化

for $t = 1, \dots, T$

$$m \leftarrow t \bmod M$$

$$U \leftarrow U + \eta(t) \nabla_U f_{l^2}(MC^{(m)}; U, V)$$

$$V \leftarrow V + \eta(t) \nabla_V f_{l^2}(MC^{(m)}; U, V)$$

end for

出力: 因子行列 U, V

欠損値がある場合の行列分解

欠損値を除外

欠損部分を完全に未知だとし、学習の際に一切使わない方法
 X の観測部分のインデックス集合 Ω

$$\text{目的関数: } \sum_{(i,j) \in \Omega} (x_{i,j} - u_i^T v_j)^2$$

$$\text{勾配: } \nabla_U f_{l^2} = -((X - UV^T) * M)V + \lambda U$$

$$\text{マスク } m_{i,j} = \begin{cases} 1 & (i,j) \in \Omega \\ 0 & \text{それ以外} \end{cases}$$

欠損値を補完

X の欠損値を何らかの形で補完した行列を \hat{X} とする

$$\text{目的関数: } \|\hat{X} - UV^T\|_{\text{Fro}}^2$$

補完の仕方はいくつかあるが、交互最適化と組み合わせる場合、
 U, V をランダムに初期化した後 $t+1$ 回目の反復において t 回目の予測値で補完

$$U^{t+1} \text{の更新: } i, j \notin \Omega \text{ において } \hat{x}_{i,j} = (u_i^{(t)})^T v_j^{(t)}$$

$$V^{t+1} \text{の更新: } i, j \notin \Omega \text{ において } \hat{x}_{i,j} = (u_i^{(t+1)})^T v_j^{(t)}$$

*「関係データ学習」 石黒勝彦・林浩平 講談社

*「続・わかりやすい パターン認識」 石井健一郎・上田修功 オーム社