

数理情報工学特論第一
【機械学習とデータマイニング】
4章：教師なし学習③

かしま ひさし
鹿島 久嗣
(数理 6 研)

kashima@mist.i.~

グラフィカルモデルについて学びます

- グラフィカルモデル
- グラフィカルラッソ
- グラフィカルラッソの推定アルゴリズム

グラフィカルモデル

教師なし学習の主要タスクは4つあるのです

- 教師なし学習においては通常、データ上の確率分布 $P(\phi(x))$ を何らかの形で推定することが行われる
- 確率分布の使い道としては主に以下の4つが挙げられる：
 - **タスク1: 確率分布そのものを用いた分析** (← 今日はコレ)
 - **タスク2: データの確率評価**
 - **タスク3: 未観測値 (欠損値) の推定** (← 前々回)
 - **タスク4: 潜在変数の推定** (← 前回)

今回は、データの特徴ベクトルの要素間の関係を調べるために、精度行列が疎であるような多次元確率分布を考えます

- 4つのタスクのひとつ「確率分布そのものを用いた分析」では、確率分布そのものを用いてデータに対する知見を得ることが目的
 1. パラメータ（もしくは確率分布 $P(\phi(x))$ の形）を「鑑賞」することによって、データ全体がどのような形で分布しているかを知る
 2. 2つのデータ集合から推定された2つの確率分布が異なるか、異なるとしたらどのように異なるのかを調べる
- 多次元正規分布モデルの場合、精度行列 Λ の各要素は、任意の2要素の積にかかる係数であるから、要素間の直接的な関係を表している
$$P(\phi(x); \mu, \Lambda) \equiv \frac{1}{(2\pi)^{D/2}} |\Lambda|^{1/2} \exp\left(-\frac{1}{2}(\phi(x) - \mu)^\top \Lambda (\phi(x) - \mu)\right)$$
- この要素間の関係は、精度行列の非零要素の数、つまり、直接的な関係の数が少ない方が、解釈が容易になる
- そのため、今回は、精度行列が疎（多くの要素が0）で、要素間の関係がよりはっきりと示されているような、多次元正規分布を扱う

一般に、パラメータがグラフ構造（疎な構造）をもつと解釈できるようなモデルをグラフィカルモデルと呼びます

- （多次元正規分布に限らず）一般に、疎なパラメータ構造をもち、その構造がある種のグラフ構造として解釈できるような確率モデルをグラフィカルモデルと呼ぶ
 - 連続の場合の代表的モデル：
 - グラフィカルラッソ（疎な多次元正規分布）
 - 離散の場合の代表的モデル：
 - ベイジアンネットワーク
- 広い意味では、教師つき学習のところで紹介した条件付き確率場（CRF）もグラフィカルモデルの範疇に入る

グラフィカルラッソ

精度行列の最尤推定量は通常「密」になります

- 多次元正規分布の精度行列を観測されたデータからの最尤推定する方法については以前述べた
- たとえ、真の精度行列が疎である場合でも、通常は推定された精度行列は密になってしまう

- 精度行列 Λ を用いた多次元正規分布の密度関数：

$$P(\phi(x); \mu, \Lambda) \equiv \frac{1}{(2\pi)^{D/2}} |\Lambda|^{\frac{1}{2}} \exp\left(-\frac{1}{2}(\phi(x) - \mu)^\top \Lambda (\phi(x) - \mu)\right)$$

- パラメータ μ および Λ を、訓練データ集合 $\{x^{(i)}\}_{i=1}^N$ から最尤推定によって求めることにすると、その最適化問題は：

$$(\mu^*, \Lambda^*) \equiv \operatorname{argmax}_{(\mu, \Lambda)} \frac{1}{N} \sum_{i=1}^N \log P(\phi(x^{(i)}); \mu, \Lambda)$$

- このままでは精度行列の推定量 Λ^* が ($=\Sigma^{*-1}$) 密行列になってしまう

$$\Sigma^* = \frac{1}{N} \sum_{i=1}^N (\phi(x^{(i)}) - \mu)(\phi(x^{(i)}) - \mu)^\top$$

精度行列の推定量を「疎」にするために、 L_1 正則化を行ったものをグラフィカルラッソと呼びます

- パラメータ μ および Λ を、訓練データ集合 $\{x^{(i)}\}_{i=1}^N$ から最尤推定によって求めることにすると、その最適化問題は：

$$(\mu^*, \Lambda^*) \equiv \operatorname{argmax}_{(\mu, \Lambda)} \frac{1}{N} \sum_{i=1}^N \log P(\phi(x^{(i)}); \mu, \Lambda)$$

– このままでは、精度行列の推定量 Λ^* が密行列になってしまう

- 精度行列の推定量 Λ^* を疎にするために、精度行列について L_1 正則化を行うことにすると、最適化問題は：

$$(\mu^*, \Lambda^*) \equiv \operatorname{argmax}_{(\mu, \Lambda)} \frac{1}{N} \sum_{i=1}^N \log P(\phi(x^{(i)}); \mu, \Lambda) - \lambda \|\Lambda\|_1$$

– $\|\Lambda\|_1$ は行列の L_1 -ノルムであり、以下のように各要素の絶対値の和

$$\|\Lambda\|_1 \equiv \sum_{k, \ell=1}^D |[\Lambda]_{k, \ell}|$$

L_1 正則化を行うと、精度行列の推定量は閉じた形では求まりません

- 平均パラメータの推定量 μ^* については正則化項が存在しないため、最尤推定の場合と同じになり、その推定量は：

$$\mu^* = \frac{1}{N} \sum_{i=1}^N \phi(x^{(i)})$$

– 全データの特徴ベクトルの平均

- 一方で、精度行列のほうは L_1 正則化項の存在により、平均パラメータのように閉じた形では求まらない
- 従って、グラフィカルラッソの本質は、いかにして精度行列の推定を行うかということになる

グラフィカルラッソの目的関数はシンプルに書けます

- Λ についてのみ最適化することにして、多次元正規分布の密度関数の定義を代入して書き下すと：

$$\Lambda^* \equiv \operatorname{argmax}_{\Lambda} \frac{1}{2} \log |\Lambda| - \frac{1}{2N} \sum_{i=1}^N (\phi(x) - \mu)^\top \Lambda (\phi(x) - \mu) - \lambda \|\Lambda\|_1$$

- この最適化問題の目的関数は以下のように書き換えられる：

$$J(\Lambda) \equiv \frac{1}{2} \log |\Lambda| - \frac{1}{2} \operatorname{Tr} \Lambda \mathbf{S} - \lambda \|\Lambda\|_1$$

— ここで、以下を使った：

- $(\phi(x) - \mu)^\top \Lambda (\phi(x) - \mu) = \operatorname{Tr}(\phi(x) - \mu)^\top \Lambda (\phi(x) - \mu)$
 $= \operatorname{Tr} \Lambda (\phi(x) - \mu) (\phi(x) - \mu)^\top$
- 標本分散共分散行列 \mathbf{S} ：

$$\mathbf{S} \equiv \frac{1}{N} \sum_{i=1}^N (\phi(x) - \mu) (\phi(x) - \mu)^\top$$

グラフィカルラッソのアルゴリズム

まずは目的関数を精度行列で微分します

- グラフィカルラッソの目的関数を最大化する Λ を求めたい

$$J(\Lambda) \equiv \frac{1}{2} \log |\Lambda| - \frac{1}{2} \text{Tr} \Lambda \mathbf{S} - \lambda \|\Lambda\|_1$$

- 実は、この最適化問題は精度行列 Λ についてではなく、その逆行列である分散共分散行列 $\Sigma = \Lambda^{-1}$ について解くことになる

- この目的関数を Λ の各 (k,l) -要素 $[\Lambda]_{k,l}$ で微分すると：

$$\frac{\partial J(\Lambda)}{\partial [\Lambda]_{k,l}} = \begin{cases} \frac{1}{2} [\Sigma - \mathbf{S}]_{k,l} - \lambda & (\text{if } [\Lambda]_{k,l} > 0) \\ \frac{1}{2} [\Sigma - \mathbf{S}]_{k,l} + \lambda & (\text{if } [\Lambda]_{k,l} < 0) \\ \text{undefined} & (\text{if } [\Lambda]_{k,l} = 0) \end{cases}$$

- ここで行列の微分の公式をもちいた：

$$\frac{\partial \log |\Lambda|}{\partial \Lambda} = \Lambda^{-1} = \Sigma \qquad \frac{\partial \text{Tr} \Lambda \mathbf{S}}{\partial \Lambda} = \mathbf{S}$$

- Λ と Σ の逆行列を通じた等価性から、↑の式において (Λ についての微分) = 0 の成立は、対応する Σ が最適解であることも意味する

対角成分は閉じた形で求められます

- まず、微分を Λ の対角成分についてみる
- 精度行列 Λ は正定行列であることに注意すると、その対角成分は必ず正、つまり $[\Lambda]_{k,k} > 0$ であることから、微分の対角成分は：

$$\frac{\partial J(\Lambda)}{\partial [\Lambda]_{k,k}} = \frac{1}{2} [\Sigma - \mathbf{S}]_{k,k} - \lambda$$

- これを0と置くと、分散共分散行列 Σ の対角成分を閉じた形で得る：

$$[\Sigma]_{k,k} = [\mathbf{S}]_{k,k} + 2\lambda$$

- これで分散共分散行列 Σ の対角成分については最適解が求まった
- 今後は Σ の非対角成分についてのみの最適化を考えればよい

非対角成分は、閉じた形で求まりませんが、実は、 L_1 正則化回帰（の繰り返し）に帰着されます

- 分散共分散行列 Σ の対角成分以外の成分の最適解を求める
 - 精度行列 Λ の最適解は $\Lambda = \Sigma^{-1}$ の関係によって自動的に定まる
- しかし、対角成分全てについての最適化は難しい
- そこで、分散共分散行列のある行（対称なので、ある列としても同じ）のみについての最適化を行い、これを選ぶ行を変えながら繰り返すという戦略をとることにする
- これから示すように、実は、分散共分散行列のある行（列）についての最適化は、回帰問題のときに紹介した L_1 正則化回帰に帰着されることがわかる

分散共分散行列を部分的に最適化するため、各行列の分割を考えます

- 分散共分散行列を部分的に最適化するため、分散共分散行列の分割を考える：

$$\Sigma = \begin{bmatrix} \tilde{\Sigma} & \tilde{\sigma} \\ \tilde{\sigma}^\top & \sigma_D \end{bmatrix}$$

- $\tilde{\Sigma}$ は $(D-1) \times (D-1)$ の行列、 $\tilde{\sigma}$ は長さ $D-1$ のベクトル、 σ_D はスカラー
- Σ の最後の行（列）について最適化するとすれば、 $\tilde{\Sigma}$ は定数であり、 σ_D は対角成分なので既に最適解が求まっており、定数とみなせる
- 従って、ここでパラメータは $\tilde{\sigma}$ となる
- Σ と同様に、精度行列と標本分散共分散行列の分割をそれぞれ：

$$\Lambda = \begin{bmatrix} \tilde{\Lambda} & \tilde{\lambda} \\ \tilde{\lambda}^\top & \lambda_D \end{bmatrix} \quad \mathbf{S} = \begin{bmatrix} \tilde{\mathbf{S}} & \tilde{\mathbf{s}} \\ \tilde{\mathbf{s}}^\top & s_D \end{bmatrix}$$

分散共分散行列と精度行列の関係式が ブロックごとに導かれます

- $\Sigma \Lambda = \mathbf{I}$ であることから、分割した行列の各ブロックについて：

$$\begin{bmatrix} \tilde{\Sigma} & \tilde{\sigma} \\ \tilde{\sigma}^\top & \sigma_D \end{bmatrix} \begin{bmatrix} \tilde{\Lambda} & \tilde{\lambda} \\ \tilde{\lambda}^\top & \lambda_D \end{bmatrix} = \begin{bmatrix} \tilde{\Sigma}\tilde{\Lambda} + \tilde{\sigma}\tilde{\lambda}^\top & \tilde{\Sigma}\tilde{\lambda} + \lambda_D\tilde{\sigma} \\ \tilde{\sigma}^\top\tilde{\Lambda} + \sigma_D\tilde{\lambda}^\top & \tilde{\sigma}^\top\tilde{\lambda} + \sigma_D\lambda_D \end{bmatrix} = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0}^\top & 1 \end{bmatrix}$$

- 特に、右上のブロックから： $\tilde{\Sigma}\tilde{\lambda} + \lambda_D\tilde{\sigma} = 0$
- すなわち： $\tilde{\lambda} = -\lambda_D\tilde{\Sigma}^{-1}\tilde{\sigma}$

$$\Sigma = \begin{bmatrix} \tilde{\Sigma} & \tilde{\sigma} \\ \tilde{\sigma}^\top & \sigma_D \end{bmatrix} \quad \Lambda = \begin{bmatrix} \tilde{\Lambda} & \tilde{\lambda} \\ \tilde{\lambda}^\top & \lambda_D \end{bmatrix}$$

分散共分散行列の現在注目しているブロックが最適解になる条件を導きます

- 先の微分の $\tilde{\lambda}$ (精度行列の現在注目しているブロック) に関連する部分だけを取り出すと：

$$\frac{\partial J(\Lambda)}{\partial [\tilde{\lambda}]_k} = \begin{cases} \frac{1}{2}[\tilde{\sigma} - \tilde{s}]_k - \lambda & (\text{if } [\tilde{\lambda}]_k > 0) \\ \frac{1}{2}[\tilde{\sigma} - \tilde{s}]_k + \lambda & (\text{if } [\tilde{\lambda}]_k < 0) \\ \text{undefined} & (\text{if } [\tilde{\lambda}]_k = 0) \end{cases}$$

- となっているので、ここから $\tilde{\lambda}$ を $\tilde{\lambda} = -\lambda_D \tilde{\Sigma}^{-1} \tilde{\sigma}$ (前頁) を用いて消去すれば：

$$\frac{\partial J(\Lambda)}{\partial [\tilde{\lambda}]_k} = \begin{cases} \frac{1}{2}[\tilde{\sigma} - \tilde{s}]_k - \lambda & (\text{if } [-\lambda_D \tilde{\Sigma}^{-1} \tilde{\sigma}]_k > 0) \\ \frac{1}{2}[\tilde{\sigma} - \tilde{s}]_k + \lambda & (\text{if } [-\lambda_D \tilde{\Sigma}^{-1} \tilde{\sigma}]_k < 0) \\ \text{undefined} & (\text{if } [-\lambda_D \tilde{\Sigma}^{-1} \tilde{\sigma}]_k = 0) \end{cases}$$

- 結果、微分が、 $\tilde{\sigma}$ (と λ_D) によって表される
 - これが =0 になるのが最適解の条件
 - しかし、このままでは、場合分けの条件式がややこしい...

パラメータを置き換えて、最適解の条件を簡単に書きます

■ この微分を簡単にしたい：
$$\frac{\partial J(\Lambda)}{\partial [\tilde{\lambda}]_k} = \begin{cases} \frac{1}{2}[\tilde{\sigma} - \tilde{s}]_k - \lambda & (\text{if } [-\lambda_D \tilde{\Sigma}^{-1} \tilde{\sigma}]_k > 0) \\ \frac{1}{2}[\tilde{\sigma} - \tilde{s}]_k + \lambda & (\text{if } [-\lambda_D \tilde{\Sigma}^{-1} \tilde{\sigma}]_k < 0) \\ \text{undefined} & (\text{if } [-\lambda_D \tilde{\Sigma}^{-1} \tilde{\sigma}]_k = 0) \end{cases}$$

■ 新たなパラメータ β を導入する： $\beta \equiv \tilde{\Sigma}^{-1} \tilde{\sigma}$

– つまり： $\tilde{\sigma} = \tilde{\Sigma} \beta$

■ であること、また、精度行列 Λ は正定であるので、その対角成分である λ_D は正であることに注意すると、微分の式は：

$$\frac{\partial J(\Lambda)}{\partial [\tilde{\lambda}]_k} = \begin{cases} \frac{1}{2}[\tilde{\Sigma} \beta - \tilde{s}]_k - \lambda & (\text{if } [\beta]_k < 0) \\ \frac{1}{2}[\tilde{\Sigma} \beta - \tilde{s}]_k + \lambda & (\text{if } [\beta]_k > 0) \\ \text{undefined} & (\text{if } [\beta]_k = 0) \end{cases}$$

– β のみを使って書くことができる。

– 場合分けの不等式の向きが変わったことに注意する

同じ微分をもつ、別の最適化問題に書き換えます

- 最終的に得られた微分：

$$\frac{\partial J(\Lambda)}{\partial [\tilde{\lambda}]_k} = \begin{cases} \frac{1}{2}[\tilde{\Sigma}\beta - \tilde{s}]_k - \lambda & (\text{if } [\beta]_k < 0) \\ \frac{1}{2}[\tilde{\Sigma}\beta - \tilde{s}]_k + \lambda & (\text{if } [\beta]_k > 0) \\ \text{undefined} & (\text{if } [\beta]_k = 0) \end{cases}$$

- これは、以下の目的関数の β についての微分と一致する：

$$\tilde{J}(\beta) \equiv \frac{1}{4}\beta^\top \tilde{\Sigma}\beta - \frac{1}{2}\tilde{s}^\top \beta + \lambda \|\beta\|_1$$

- もともとの最大化問題の目的関数を β の関数としてみたものと、この最小化問題の目的関数は同一
 - 連続で、微分不可能な点が同じで、それ以外の場所では微分が一致
- つまり、もともとの目的関数のかわりに、新たな目的関数を用いて β を求めても差し支えないということの意味する

新しい目的関数を、変数ごとに最適化することを考えます

- さらに、目的関数：

$$\tilde{J}(\beta) \equiv \frac{1}{4}\beta^\top \tilde{\Sigma}\beta - \frac{1}{2}\tilde{s}^\top \beta + \lambda \|\beta\|_1$$

を、 β の第 k 成分 $[\beta]_k$ についてのみ最適化することを考える

- 目的関数 $\tilde{J}(\beta)$ を $[\beta]_k$ についての関数だと思って整理すると：

$$\begin{aligned}\tilde{J}([\beta]_k) &\equiv \frac{1}{4}[\tilde{\Sigma}]_{k,k}[\beta]_k^2 + \frac{1}{2} \sum_{j \neq k} [\beta]_j [\tilde{\Sigma}]_{j,k} [\beta]_k - \frac{1}{2}[\tilde{s}]_k [\beta]_k + \lambda |[\beta]_k| + \text{const.} \\ &= \frac{1}{4}[\tilde{\Sigma}]_{k,k}[\beta]_k^2 - \frac{1}{2} \left([\tilde{s}]_k - \sum_{j \neq k} [\beta]_j [\tilde{\Sigma}]_{j,k} \right) [\beta]_k + \lambda |[\beta]_k| + \text{const.} \\ &= \frac{1}{4}[\tilde{\Sigma}]_{k,k} \left([\beta]_k - \frac{[\tilde{s}]_k - \sum_{j \neq k} [\beta]_j [\tilde{\Sigma}]_{j,k}}{[\tilde{\Sigma}]_{k,k}} \right)^2 + \lambda |[\beta]_k| + \text{const.}\end{aligned}$$

ちょうどL₁正則化回帰のときに出てきた目的関数と同じ形であることに気付きます

- 結局、 $[\beta]_k$ の最適解 $[\beta]_k^*$ は、以下の最適化問題を解くことで求まる

$$[\beta]_k^* \equiv \operatorname{argmin}_{[\beta]_k} \frac{1}{2} ([\beta]_k - [\tilde{\beta}]_k)^2 + \gamma |[\beta]_k|$$

$$\text{where } \gamma \equiv \frac{2\lambda}{[\tilde{\Sigma}]_{k,k}} \text{ and } [\tilde{\beta}]_k \equiv \frac{[\tilde{s}]_k - \sum_{j \neq k} [\beta]_j [\tilde{\Sigma}]_{j,k}}{[\tilde{\Sigma}]_{k,k}}$$

- これはまさにL₁正則化回帰の次元ごとの逐次解法のとときに出てきた最適化問題↓と同じ形をしている

$$w_d^* = \operatorname{argmin}_{w_d} \tilde{L}(w_d)$$

$$\tilde{L}(w_d) \equiv \frac{1}{2} (w_d - \tilde{w}_d)^2 + \gamma |w_d|$$

L₁正則化回帰のときの解を利用して、解を導きます

- L₁正則化回帰の次元ごとの逐次解法の際に出てきた最適化問題：

$$w_d^* = \operatorname{argmin}_{w_d} \tilde{L}(w_d) \text{ where } \tilde{L}(w_d) \equiv \frac{1}{2} (w_d - \tilde{w}_d)^2 + \gamma |w_d|$$

- この解は：
$$w_d^* = \begin{cases} \tilde{w}_d - \gamma & (\text{if } \tilde{w}_d > \gamma) \\ \tilde{w}_d + \gamma & (\text{if } \tilde{w}_d < -\gamma) \\ 0 & (\text{if } -\gamma \leq \tilde{w}_d \leq \gamma) \end{cases}$$

- これを利用すると、いま解きたい最適化問題：

$$[\beta]_k^* \equiv \operatorname{argmin}_{[\beta]_k} \frac{1}{2} ([\beta]_k - [\tilde{\beta}]_k)^2 + \gamma |[\beta]_k|$$

の最適解 $[\beta]_k^*$ は：

$$[\beta]_k^* = \begin{cases} [\tilde{\beta}]_k - \gamma & (\text{if } [\tilde{\beta}]_k > \gamma) \\ [\tilde{\beta}]_k + \gamma & (\text{if } [\tilde{\beta}]_k < -\gamma) \\ 0 & \text{if } (-\gamma \leq [\tilde{\beta}]_k \leq \gamma) \end{cases}$$

- これを k を変えながら繰り返すことによって、 β が求まる

L_1 正則化回帰を繰り返すことによって、グラフィカルラッソの推定ができます

- 一旦 β が得られると、 $\tilde{\sigma} = \tilde{\Sigma}\beta$ を用いて、 $\tilde{\sigma}$ を求めることができる
- この一連の手続きによって Σ の最後の行（列）を求めるのがグラフィカルラッソのアルゴリズムにおける1ステップである
- このステップを、行（列）の順番を入れ替えて最後にくるものを適当に変えながら繰り返し行い、収束するまで繰り返す
 - つまり、 L_1 正則化回帰問題を繰り返し解く
- グラフィカルラッソのアルゴリズムは：
 - ブロック分割で最後にもってくる行（列）を変えながら、最後の行（列）についての最適化を繰り返す外側のループ
 - 最後の行（列）についての最適化を行うために、 β のある次元についての最適化を、選ぶ次元を変えながら繰り返す内側のループ

の2つのループで構成されている

このアルゴリズムは逆行列の計算を必要としないため効率的です

- このアルゴリズムの効率的なところは、 $\Sigma = \Lambda^{-1}$ という関係をもつ2つの行列 Σ と Λ を扱っているながら、逆行列の計算がまったく出てこないところである
 - $\Sigma \Lambda = \mathbf{I}$ をブロック分割した式（右上と右下）および $\tilde{\sigma} = \tilde{\Sigma} \beta$ の連立方程式を解くことで $\tilde{\sigma}$ から Λ の最後の行（列）が求まる：

$$\lambda_D = \frac{1}{\sigma_D - \tilde{\sigma}^\top \beta} \quad \tilde{\lambda} = -\frac{\beta}{\sigma_D - \tilde{\sigma}^\top \beta}$$

- また、得られた解は条件 $\Sigma \Lambda = \mathbf{I}$ を満たすように作られるため、得られる Σ と Λ は正定であることが保証される