スパース表現の数理とその応用

日野 英逸^{1,a)} 村田 昇^{1,b)}

概要:スパースコーディングは生物の一次視覚野の情報処理を数学的にモデル化したものであり,与えられた画像を少数の基底の線型結合で表現する手法である.観測信号のスパース表現は,工学的にも効率的な情報の保持・伝達,あるいはノイズに対して頑健な情報表現を実現する手法として注目を集めている. 本稿では,スパースコーディングを始めとする種々の行列分解手法の数理的側面を,その確率モデルを介して統一的に論じる.また,スパースコーディングの代表的なアルゴリズムと幾つかの応用を紹介する.

1. はじめに

生物の一次視覚野の細胞は,方向と空間スケールに関し て選択的な空間フィルタを有していることが古くから知ら れていた. 一次視覚野における受容細胞は, ある方向を向 いたある空間周波数成分が網膜上の特定の領域に出現する と,選択的に反応する性質を持っている.細胞が反応する 方向と周波数は受容細胞毎に異なるため,視覚野全体とし て空間フィルタバンクを形成していると考えることが出 来る.一次視覚野の細胞がこうした特徴を持つに至った理 由を説明するために様々な試みがなされてきたが,ここ十 数年の間に,自然画像の統計的構造を積極的に利用するこ とによって自然画像を効率的に符号化 (コーディング) す るための仕組みとして獲得されたとする考え方が提案さ れ脚光を浴びている [1,2]. Olshausen らによる一連の研 究では、「自然画像を基底画像の線型結合で表した時、そ の結合係数が疎である」という制約の下で基底を学習する ことで,生物の一次視覚野における受容細胞と同様の特徴 を有する基底が得られることが実験的に示されている.本 稿では画像の小領域の画素値を観測信号として扱い,信号 のスパース表現の問題を解説する.例えば一辺が \sqrt{d} ピク セルの画像の小領域を扱う場合には, $\sqrt{d} imes \sqrt{d}$ の画素値 をベクトル表現した $x \in \mathbb{R}^d$ が一つの観測信号ということ になる.Olshausen & Field [1] は,画像 x が基底ベクト $\boldsymbol{\nu}^{*1}\boldsymbol{d}_k, k = 1, \dots, m$ の線型結合

$$\boldsymbol{x} = \sum_{k=1}^{m} c_k \boldsymbol{d}_k, \ c_k \in \mathbb{R}$$
(1)

□ 早稲田大学理工学術院 〒 169-8555 東京都新宿区大久保 3-4-1

- $^{\rm a)} \quad {\rm hideitsu.hino}@toki.waseda.jp$
- ^{b)} noboru.murata@eb.waseda.ac.jp
- *1 必ずしも一次独立なベクトルとは限らず,基底という表現は正確ではないが,慣習上「基底」と呼ぶことにする.

© 2012 Information Processing Society of Japan

で表されるという仮定を置いた.この基底系 $\{d_k\}$ は,画 像を「効率的に」表現することが出来るように構成する. ここでいう効率的とは,ある画像を表現する際に係数ベク トル $c = (c_1, ..., c_m) \in \mathbb{R}^m$ のうちの少数の c_k のみが非ゼ 口の値を取り,残りの大部分はゼロとなることを意味する. もし非ゼロの c_k の数がdより少なく,かつ基底系 $\{d_k\}$ を 予めうまく選ぶことができているとすれば,d 個の値を持 つ画像xをより少ない個数の非ゼロの c_k の集合で効率良 く表現することができるからである.このように,意味の ある非ゼロ要素が全体に対して少数である状態を,スパー ス(sparse; 疎)であると呼ぶ.適切な基底系 $\{d_k\}$ とその結 合係数 $\{c_k\}$ を得るには,観測信号の情報を保存する項(文 献 [1]では l_2 ノルムを考えている)と,結合係数がスパー スになることを促進する正則化項 Ψ の和

$$\min_{\{\boldsymbol{d}_k\},\{c_k\}} \|\boldsymbol{x} - \sum_{k=1}^m c_k \boldsymbol{d}_k\|_2^2 + \Psi(\boldsymbol{c})$$
(2)

の最小化を行えば良い.例えば [1] では Ψ として $\Psi(c) = \lambda \sum_{k=1}^{d} e^{-c_k^2}, \lambda > 0$ を用いることを提案している.こうして得られた基底は,生物の一次視覚野細胞とよく似た局所的な空間周波数特性を持つ Gabor wavelet 基底と類似したものになる [3].

なお,Olshausen & Field [1] 以前から,ニューラルネットワークや連想記憶の研究において,スパースな符号化に よって情報の表現及び保持を効率的に行うことが出来るこ とが知られていた.例えば [4-6] では,工学的観点からス パース表現の利点の定量的議論が行われている.

以下,本稿の構成を記す.第2節では画像の生成モデル を簡単に説明する.第3節ではスパースコーディングの定 式化を行ない,スパース表現において中心的役割を果たす ノルムを導入する.第4節ではスパース表現を観測信号行 列の分解手法として捉え,関連する種々の手法を概説し, 第5節でこれらの行列分解手法の背後にある確率モデルを 説明する.第6節及び第7節ではそれぞれスパース表現 のための係数,基底の最適化手法を紹介し,第8節ではス パース表現の画像処理における応用例を挙げる.

2. 画像の生成モデル

本稿ではスパース表現の具体的な例として画像を取り上 げる.その準備のために,画像の生成モデルとして重要な 2つの考え方を簡単に纏める.

2.1 ピクセルベースモデル

多数の変数とその変数間の相互作用からなる系のうち, 近傍のみで条件付けられる単純なモデルの一つとしてはマ ルコフ確率場 (Markov Random Field; MRF) があり,画 像の確率モデルとしても広く利用されている[7].このモ デルは,あるピクセルのとる値vの確率がその周辺のピク セル値で決まる条件付き確率

p(v|周辺のピクセル値)

によって表されると考える,ピクセルベースのモデルである.マルコフ確率場は事前分布として画像の滑らかさなど に対応する制約を容易に取り込むことができる.例えば劣 化画像を観測した場合に,観測した画像で条件付けた原画 像の生成確率(事後確率)を利用して画像復元を行う方法 や,超解像への応用が知られている[8,9].また,局所的な 相互作用を記述した確率モデルを導入することで,統計力 学の分野で発展した平均場近似を始めとする各種の近似計 算手法が利用できる[8].しかし,実用に供するレベルで の画像の生成機構をモデル化するためには,事前分布の適 切な設定と,本質的に非凸な問題の最適化技術が必要とな り,綿密な設計と高度なノウハウが要求される.

2.2 基底ベースモデル

画像を基底画像の線型結合により表現するアプローチ を,以下では基底ベースモデルと呼ぶ.画像全域を表現す る基底を利用することもあるが,一般には画像の小領域を 対象としてモデル化を行い,小領域をずらしつつ画像全域 を表現することが多い.特に自然画像はある程度小さな領 域に限ると類似したパターンを示すことが多く,従って

x = Dc, # $\{i | c_i \neq 0\} \ll d (= 領域の画素数)$,

のように少数の基本的な基底の組合せで表現できる.ここ で #S は集合 S の要素数を表す.また,多くの画像処理手 法の計算量は観測信号の次元に対して指数関数的に増加す るため,小領域毎に処理を行う基底ベースモデルには計算 量の低減という利点もある.本稿ではスパース表現の数理 的側面の説明を,基底ベースモデルによる画像の表現を前 提として行う.

3. スパースコーディングの定式化

本節ではスパース表現の代表例としてスパースコーディング問題を定式化する.また,スパース性と関連の深いベクトル及び行列のノルムについて説明する.

3.1 定式化

 $d 次元観測信号ベクトルを <math>x \in \mathbb{R}^d$ とし, m 個の基底を並 べた辞書 $D = (d_1, \dots, d_m)$ の線型結合により x = Dcの ように表現することを考える.この時,スパースコーディ ングの枠組みでは $m \gg d$ という過剰決定系^{*2}による近似 を考える.信号の次元より多い基底による表現 x = Dc で は c の一意性を保証することが出来ないので,通常は観測 信号 x の表現に利用される基底を D のうちの一部に制限 する.つまり, $\|c\|_0$ で c の l_0 ノルム, すなわちベクトル c の非ゼロ成分の数を表すとして,スパースコーディング は典型的には最適化問題

$$\min_{\boldsymbol{c}} \|\boldsymbol{x} - D\boldsymbol{c}\|_{2}^{2} + \lambda \|\boldsymbol{c}\|_{0}, \quad \lambda > 0$$
(3)

として定式化される.しかしながら,この問題は全ての基 底の組み合わせを試さないと最適解が得られない組合せ最 適化問題であり,NP困難であることが知られている[10]. そこで,l₁ノルムへの緩和問題

$$\min_{\boldsymbol{c}} \quad \|\boldsymbol{x} - D\boldsymbol{c}\|_2^2 + \lambda \|\boldsymbol{c}\|_1, \quad \lambda > 0$$
(4)

を考えることが多い.このl₁ノルム正則化問題は線型計画 問題として表現することが可能である.

より一般に,辞書と係数を求める問題は統計的推測の 枠組みで次のように表される.n 個の観測信号を並べた行 列を $X = (x_1, \cdots, x_n)$,辞書行列を $D = (d_1, \ldots, d_m) =$ $(d^1, \ldots, d^d)^\top$,各信号を表現するための辞書の係数をな らべた係数行列を $C = (c_1, \cdots, c_n) = (c^1, \cdots, c^m)^\top$ とす る.ここで, $d_i \in \mathbb{R}^d$ は行列Dの第i列, $d^j \in \mathbb{R}^m$ は行 列Dの第j行を表し,同様に $c_i \in \mathbb{R}^m$ は行列Cの第i列, $c^j \in \mathbb{R}^n$ は第j行を表すものとする.L(D,C;X)で近似の 損失関数を表し,信号の表現がスパースならば小さい値を 取る正則化項 $\Psi(D,C)$ を用いて,n個の観測信号に対する スパースコーディングは最適化問題

$$\min_{D,C} \quad L(D,C;X) + \Psi(D,C) \tag{5}$$

として定式化される.後述するように,損失関数は負の対 数尤度, Ψ は基底あるいは係数に対する事前分布の負の対 数尤度として確率モデルでの解釈が可能であることが多い. 本稿で主に扱う式(5)の定式化は,損失関数L(D,C;X)を 可能な限り小さくすることで信号の再構成を精度よく行い つつ,正則化項 $\Psi(D,C)$ の最小化により表現のスパース性 ^{*2} 機械学習の文脈では過完備基底 (overcomplete basis) とも呼ぶ. を確保しようというものである.問題を実際に解く際や特定の理論解析においては,例えば

$$\begin{array}{ll} \min_{D,C} \quad \Psi(D,C) \quad \text{subject to} \quad L(D,C;X) < \epsilon, \\ \min_{D,C} \quad L(D,C;X) \quad \text{subject to} \quad \Psi(D,C) < \epsilon \end{array}$$

のような定式化が便利なこともある.種々の表現の同等性,及び用途に応じた適切な表現に関しては[11]に詳しくまとめられている.

3.2 ベクトルのノルム

基底による信号の表現でのスパース性は,基本的には結 合係数に対するノルム制約から導かれる.本稿では一般に ベクトル $x \in \mathbb{R}^d$ に対して,

$$\|\boldsymbol{x}\|_{p} = \left(\sum_{i=1}^{d} |x_{i}|^{p}\right)^{1/p} = \sqrt[p]{|x_{1}|^{p} + \dots + |x_{d}|^{p}} \quad (6)$$

で l_p ノルム (l_p norm) を表現する.特にp = 0の時は

$$\|\boldsymbol{x}\|_{0} = \lim_{p \to 0} \|\boldsymbol{x}\|_{p}^{p} = \lim_{p \to 0} \sum_{i=1}^{d} |x_{i}|^{p} = \#\{i | x_{i} \neq 0\}$$
(7)

とする.定義 (6) において $0 \le p < 1$ の範囲では厳密には $\|x\|_p$ はノルムの公理を満たさない.すなわち, p = 0の時 は斉次性

 $\|a\boldsymbol{x}\|_p = |a| \|\boldsymbol{x}\|_p, \quad a \in \mathbb{R}$

が満たされず,0<p<1の時は三角不等式(劣加法性)

$$\|m{x}_1 + m{x}_2\|_p \le \|m{x}_1\|_p + \|m{x}_2\|_p, \quad m{x}_1, m{x}_2 \in \mathbb{R}^m$$

が満たされないので, $\|x\|_p$, $0 \le p < 1$ は正確にはノルムとはならないが,慣習上 $0 \le p < \infty$ に対して $\|x\|_p \ge l_p$ ノルムと呼ぶ.

3.3 行列のノルム

行列 $M \in \mathbb{R}^{n \times m}$ に対するノルムにも種々の定義がある. 行列特有のノルムとして特に重要なものは核ノルム (nuclear norm) あるいはトレースノルム (trace norm) と呼ばれるもので,

$$\|M\|_* = \operatorname{Tr}(\sqrt{M^{\top}M}) = \sum_{i=1}^m \sigma_i \tag{8}$$

で定義される.ここで, σ_i は行列 Mの特異値である.ベクトルのノルムにおいて p = 0の場合に対応するのが,行列のランク (rank)

$$r = \#\{i|\sigma_i \neq 0\}\tag{9}$$

である.すなわち,行列 M の非ゼロ特異値の数が小さい という制約は,行列にスパースな構造を与える.しかし, ベクトルに対する l₀ ノルム制約が組み合せ最適化問題を伴 い計算が困難であるように,行列に対するランク制約は非 凸な制約となり取り扱いが難しい.そこで,ベクトルの l_0 ノルムを l_1 ノルムで緩和するように,ランク制約を核ノル ム制約で緩和することが多い.

また,行列の成分毎に定義した l_p ノルム

$$||M||_p = \left(\sum_{i=1}^n \sum_{j=1}^m |M_{ij}|^p\right)^{1/p}$$
(10)

もよく利用される.特に p = 2の場合をフロベニウスノルム (Frobenius norm) と呼ぶ.ベクトルあるいは行列のノルムに関する議論は,例えば [12] に詳しい.本稿では行列のノルムとしては核ノルムあるいは成分毎の l_p ノルムを主に利用する.

4. 行列分解としての理解

本節では,代表的なスパース表現を観測信号行列の分解 手法という観点から整理する.一般に,スパース表現は観 測信号行列 X を

$$X = DC + E,\tag{11}$$

のように辞書行列 D とその係数行列 C で表現するもので ある.ここで, $E = (\epsilon_1, \ldots, \epsilon_n), \epsilon_i \in \mathbb{R}^d$ であり誤差を表 す行列であるが,一般に各 ϵ_i は適当な平均ゼロの多変量正 規分布 $\mathcal{N}(\mathbf{0}, \sigma^2 I_d)$ に従うとすることが多い.ここで I_d は d 次元単位行列を表す.

通常は係数行列 C の各列がスパース, すなわち利用される基底が少数であるという状況を考えるが, 辞書行列 D の各列, すなわち各基底がスパースである, あるいは D, C の両方がスパースであるという問題設定もある.

4.1 スパースコーディング

スパースコーディング (Sparse Coding; SC) では多く の場合損失関数としてフロベニウスノルム $L(D,C;X) = ||X - DC||_2^2$ を考え,正則化項 $\Psi(D,C)$ としては係数ベク トル c の l_p ノルムで $0 \le p \le 1$ の場合を考える:

$$\min_{D,C} \quad \sum_{i=1}^{n} \|X - DC\|_{2}^{2} + \lambda \sum_{i=1}^{m} \|\boldsymbol{c}_{i}\|_{p}^{p}, \quad \lambda \ge 0.$$
(12)

狭義のスパースコーディングは,与えられた辞書Dを用いて信号xを近似するためのスパースな係数cを求める問題であるが,辞書の最適化も含めて考えることもある.係数C,辞書Dの具体的な最適化手法についてはそれぞれ第5節,第6節で簡単に説明する.

4.2 主成分分析

ここでは,観測信号の各次元は平均ゼロに正規化されて いるとする.すなわち,観測行列 $X = (x_1, ..., x_n)$ の各 行の和は0であると仮定する.Xのランクをrとして,X の特異値分解を

 $X = U \Sigma V^\top$

$$U = (\boldsymbol{u}_1, \dots, \boldsymbol{u}_r), \ \boldsymbol{u}_i \in \mathbb{R}^d, U^\top U = I_r$$
$$V^\top = (\boldsymbol{v}_1, \dots, \boldsymbol{v}_n), \ \boldsymbol{v}_i \in \mathbb{R}^r, V^\top V = I_r$$
$$\Sigma = \operatorname{diag}(\sigma_1, \dots, \sigma_r), \sigma_i \in \mathbb{R}$$

として,特異値は降順 $\sigma_1 \geq \sigma_2 \geq \cdots \geq \sigma_r$ に並んでいるものとする.正規直交基底系をなす辞書行列 D = Uを固定した時,係数行列は $C = \Sigma V^{\top}$ であり,信号 x_i を $x_i = Dc_i$ で表すとき,その係数ベクトルは

$$\boldsymbol{c}_i = \Sigma \boldsymbol{v}_i = (\sigma_1 v_{1i}, \dots, \sigma_r v_{ri})^{\top}$$

で与えられる.主成分分析 (Principal Component Analysis; PCA [13]) では $L(D,C;X) = ||X - DC||_2^2$ を想定し, 辞書 D は正規直交基底系をなすという制約がある.この 制約を $\Psi(D,C) = ||D^\top D - I_r||_2^2$ で表すと,主成分分析は

$$\min_{D,C} \|X - DC\|_2^2 + \lambda \|D^\top D - I_r\|_2^2$$
(13)

という行列分解問題として表現できる.辞書をD = Uとした時にこの目的関数を最小にするような係数Cは $C = \Sigma V^{\top}$ で与えられる.また,Xをランクk < rの行列の積で近似するような場合には,小さな特異値に対応する σ_i のk + 1からrまでを0にした行列 $\tilde{\Sigma} = \text{diag}(\sigma_1, \cdots, \sigma_k, 0, \cdots, 0)$ を用いてスパースな係数が

$$\boldsymbol{c}_i = \tilde{\Sigma} \boldsymbol{v}_i = (\sigma_1 v_{1i}, \sigma_2 v_{2i}, \dots, \sigma_k v_{ki}, 0, \dots, 0)^{\top}$$

で与えられる.主成分分析そのものは,観測信号行列の低 ランク行列による近似という意味でのスパース表現である が, Cの各列に対する l_p ノルム制約を考えることで積極的 にスパースな係数を得る主成分分析の拡張法も提案されて いる [14].

4.3 独立成分分析

独立成分分析 (Independent Component Analysis; ICA [15]) においては, D は辞書というよりは混合行 列であり, n 個の m 次元原信号 $C = (c_1, \ldots, c_n) =$ $(c^1, \ldots, c^m)^\top$, $c_i \in \mathbb{R}^m, c^k \in \mathbb{R}^n$ が, 混合行列 D = $(d^1, \ldots, d^d)^\top$ によって混合した結果 X が観測されると 考える.すなわち, x_k の第 i 成分は, m 次元の原信号 c_k が 混合比 $d^i = (d_1^i, \ldots, d_m^i)^\top$ で混合された信号 $(d^i)^\top c_k$ とな る.このとき, $d \ge m$, すなわち観測信号の次元が原信号の 次元以上で, 原信号 $c^k, k = 1, \ldots, m$ が非正規で互いに独 立であるならば, 観測信号 X から原信号 C を復元すること が出来る.簡単のために, $D \in \mathbb{R}^{m \times m}$ かつ正則であるとす る, すなわち観測信号の次元と原信号の次元は一致してお り, 混合行列 D は逆行列を有すると仮定する.このとき, 復元行列を $B = D^{-1}$ と置くと, $C = BX = (c^1, \dots, c^m)^\top$ であり, このCの各行 c^k が独立になるようにBを推定する.スパースな行列分解とは異なるように見えるが,次節で述べるように係数Cのエントロピー最小化の観点からスパース表現としての理解が可能である.例えば $L(D,C;X) = ||X - DC||_2^2$ とおき,確率密度関数pを持つ確率変数のエントロピーを

$$H(p) = -\int p(\boldsymbol{x}) \log p(\boldsymbol{x}) \mathrm{d}\boldsymbol{x}$$

と記すことにする. 混合行列 D の関数としての原信号 $C = D^{-1}X$ の第 j 行が従う分布を p_{c^j} とし, $\Psi(D,C) = \sum_{j=1}^m H(p_{c^j})$ を用いると, 独立成分分析の一つの定式化として次の最適化問題

$$\min_{D} \quad \|X - DC\|_{2}^{2} + \sum_{j=1}^{m} H(p_{c^{j}})$$
(14)

が得られる.

4.4 非負行列因子分解

データ解析において、データ、基底、結合係数ともに 非負に制限することが望ましいことがある。例えば、画 像のピクセル値、エネルギー、イベントの生起回数など は非負値のみを取るため、これらのデータの基本要素と なる非負の基底の加法のみで観測信号が表現される分解 が自然である。非負行列因子分解(Non-negative Matrix Factorization; NMF [16,17])は、非負データを、非負成分 のみからなる辞書と係数に分解する手法である。通常、非 負行列因子分解は正則化項 $\Psi(D,C)$ に直接非負性は含め ず、最適化問題(5)の定義域を制限して

$$\min_{[D]_{ij} \ge 0, [C]_{ij} \ge 0} \quad \|X - DC\|_2^2 + \Psi(D, C) \tag{15}$$

という問題を考える.ただし $[M]_{ij}$ は行列 M の (i, j) 成分 である.損失関数 L(D, C; X) として例えば [17] では,

$$L(D,C;X) = ||X - DC||_2^2$$
(16)

あるいは

$$L(D,C;X) = \sum_{i=1}^{d} \sum_{j=1}^{n} \left([X]_{ij} \log \frac{[X]_{ij}}{[DC]_{ij}} - [X]_{ij} + [DC]_{ij} \right) \quad (17)$$

を用いて,定義域に関する制約のみで正則化項を含まない 単純な最適化アルゴリズムを導出している.これらの損失 関数は,それぞれ,近似誤差の分布に正規分布を仮定した 場合の負の対数尤度と,平均 [*DC*]_{*ij*}のポアソン分布を仮 定した場合の一般化 Kullback-Leibler (KL) divergence に 対応する.

非負行列因子分解では非負成分のみを有する基底の加

法のみで信号を表現するため,局所性を有する辞書と,ス パースな係数が得られやすい.ただし,問題によっては解 釈が容易なスパース表現が得られにくいため,DとCに 同時にl₁ノルム正則化を施し,例えば

$$\Psi(D,C) = \lambda_1 \sum_{j=1}^m \|\boldsymbol{d}_j\|_1 + \lambda_2 \sum_{i=1}^n \|\boldsymbol{c}_i\|_1, \ \lambda_1, \lambda_2 > 0 \ (18)$$

とすることで,積極的にスパースな表現を得る手法も提案されている [18].また,L(D,C;X)として通常の KL divergence に替えて α -divergence や Bregman divergence を用いた一般化も行われている [19,20].

5. 確率モデルとしての理解

本節では,前節で説明した各種のスパース行列分解の確 率モデルとしての解釈を与える.なお,各手法は異なる複 数の確率モデルとしての解釈が可能であり,本節で示す解 釈はその一つであることに注意されたい.

5.1 スパースコーディング

簡単のため,まず辞書Dは与えられている状況を考える.損失関数を $L(D,C;X) = ||X - DC||_2^2$ としたとき,これは各信号 $x_i, i = 1, \dots, n$ の近似誤差 ϵ_i に正規分布を仮定していることになる.すなわち $p(x_i - Dc_i) = p(x_i|c_i) \propto \exp(-\frac{1}{\sigma^2} ||x_i - Dc_i||_2^2)$ である.係数の分布としては,0 の時は一般化ガウス分布

$$p(\boldsymbol{c}_i) \propto \exp\left(-\frac{\lambda}{\sigma^2} \|\boldsymbol{c}_i\|_p^p\right)$$
 (19)

を仮定していると考えられる.0<p<2のときに,この 分布は正規分布と比較して裾の重く尖度が高い優ガウス的 分布となり,スパースな解を与える*3.これにより,観測 信号と基底の結合係数の負の対数尤度は,定数項を除いて

$$-\sum_{i=1}^{n} \log p(\boldsymbol{x}_{i} | \boldsymbol{c}_{i}) p(\boldsymbol{c}_{i}) = \|X - DC\|_{2}^{2} + \lambda \sum_{i=1}^{n} \|\boldsymbol{c}_{i}\|_{p}^{p} (20)$$

となり,これを C に関して最小化することは結合係数の MAP 推定と等価である.p = 0の時に確率的解釈を行う ためには非ゼロの値を取る係数の数,係数の添字,係数の 値に対する分布を定める必要があり,特に係数の添字に関 する分布の評価が困難となる.Bayes 統計の文脈で変数選 択に利用されてきた Spike & Slab モデル [21] が,係数に l_0 正則化を課したスパースコーディングの確率モデルとし て注目を集めている (図 1).これは,係数ベクトル c の分 布として適当な連続分布 (正規分布を用いることが多い) と,中心をゼロとするディラックのデルタ関数 $\delta_0(\cdot)$ の混 合分布

*³ 典型的には, p = 1の時にラプラス分布になる.



図 1 l_0 正則化に対応する Spike & Slab 事前分布.



$$p(c_i | \sigma^2, q) = q \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}c_i^2\right) + (1-q)\delta_0(c_i) \quad (21)$$

を用いるものである.このモデルでは非ゼロのパラメタの 数を混合比 q でコントロールすることが可能であり, l₀ ノ ルムを模していると解釈できる.基底 D の事前分布にこ のモデルを仮定するもの [22] も,係数 C の事前分布に仮 定するもの [23] もある.例えば係数 C の事前分布に Spike & Slab モデルを仮定する場合は,

$$p([C]_{ij}|\sigma_c^2, q) = q \frac{1}{\sqrt{2\pi\sigma_c^2}} \exp\left(-\frac{1}{2\sigma_c^2} [C]_{ij}^2\right) + (1-q)\delta_0([C]_{ij})$$
(22)

であり,式 (5)における正則化項 $\Psi(D,C)$ としては

$$\Psi(D,C) = -\log\left\{q\frac{1}{\sqrt{2\pi\sigma_c^2}}\exp\left(-\frac{[C]_{ij}^2}{2\sigma_c^2}\right) + (1-q)\delta([C]_{ij})\right\}$$
(23)

が考えられる.

5.2 主成分分析

主成分分析は近似誤差として正規分布を仮定しているが, 辞書 D 及び係数行列 C としては直交性の制約のみを考え ており,陽に D,C の分布を考えたモデルではない.その ため式(5)に則って主成分分析の確率モデルを表現するこ とは困難であるが,以下のようにして確率モデルとして解 釈することが可能である[24]^{*4}.観測信号は平均がゼロに ^{*4} この表現は因子分析(Factor Analysis; FA [25])と関連が深い.

正規化されていると仮定して,信号の表現

$$\boldsymbol{x} = D\boldsymbol{c} + \boldsymbol{\epsilon} \tag{24}$$

における誤差 ϵ は正規分布 $\mathcal{N}(\mathbf{0}, \sigma^2 I_d)$ に従い, さらに c が 正規分布 $\mathcal{N}(\mathbf{0}, I_m)$ に従うと仮定する.このとき, $p(\mathbf{x}|\mathbf{c})$ は $\mathcal{N}(Dc, \sigma^2 I_d)$ であり, $p(\mathbf{x}) = \int p(\mathbf{x}|\mathbf{c})p(\mathbf{c})d\mathbf{c}$ は共分散 行列 $\Sigma = DD^\top + \sigma^2 I_d$ の正規分布 $\mathcal{N}(\mathbf{0}, \Sigma)$ に従うことが 分かる.このモデルを,確率的主成分分析 (Probabilistic PCA; PPCA) と呼ぶ.確率的主成分分析では,観測信号 $x_i, i = 1, \dots, n$ の尤度を最大化するように D 及び σ^2 を定 めることが出来る.すなわち,辞書 D 及び残差の分散 σ^2 をパラメタとするモデル $p(\mathbf{x})$ に観測信号行列 X を代入し て得られる負の対数尤度

$$-\log p(X|D,\sigma^2) = 2n\log|\Sigma| + \operatorname{Tr}(\Sigma^{-1}XX^{\top}), \quad (25)$$

ただし

$$\Sigma = DD^{\top} + \sigma^2 I_d$$

を D, σ^2 に関して最小化すればよい.ただし,対数尤度に おいて最適化に寄与しない定数項は省略した.

確率的主成分分析には,欠損データの取り扱い,モデル 選択,混合モデルへの拡張などが自然に行えるという利点 がある.なお,[14]に対応して,特に係数にスパース性の 制約を加えた確率的主成分分析も提案されている[26].

5.3 独立成分分析

中心極限定理によれば,有界な分散を持つ独立な確率変数の和の分布は,足し合わせる変数の数の増加とともに正規分布に近づいていく.一方,同じ分散を持つ連続確率変数の中で最大のエントロピーを持つ分布は正規分布である[27].これらの事実から,原信号が重なりあうとその分布は正規分布に近づき,エントロピーは増大することが分かる.従って,直観的には観測信号 x から推定される復元信号 Bx の各成分のエントロピーを減少させるように復元行列 $B = D^{-1}$ を選べば良いことが分かる.また,ラプラス分布を代表とする一般化ガウス分布(19)など,平均ゼロの尖度の高い優ガウス的分布はゼロ周辺に高い確率密度を有することから,スパースな信号となる.

今,観測信号 $x \in \mathbb{R}^d$ の従う確率分布の密度関数 p は 推定出来ているとする.このとき,観測信号 x と原信号 $c \in \mathbb{R}^m$ の関係から, c の確率分布は復元行列 B と確率密 度関数 p で表される.これを p_c と書くことにする.また, 復元信号 c の各成分 c^j の周辺分布を p_{c^j} と書くことにする と,原信号の各成分が独立であるとは,同時分布の密度関 数と周辺分布の密度関数の積が等しいこと,すなわち

$$p_c = p_{c^1} \cdots p_{c^n}$$

が成り立つことである.ここで,密度関数 *p*,*q* で表される 2つの分布の「離れ具合」の尺度として KL divergence [27]

$$KL(p|q) = \int p(\boldsymbol{x}) \log \frac{p(\boldsymbol{x})}{q(\boldsymbol{x})} \mathrm{d}\boldsymbol{x}$$

を採用すると,独立性の必要十分条件は同時分布と周辺分 布のエントロピーを用いて

$$KL(p_{c}||\prod_{j=1}^{m} p_{c^{j}}) = \sum_{j=1}^{m} H(p_{c^{j}}) - H(p_{c}) = 0$$

と表すことが出来る.同時分布とその周辺分布のエントロ ピーの間には不等式

$$H(p_{\boldsymbol{c}}) \leq \sum_{j=1}^m H(p_{c^j})$$

が成立するため,周辺分布のエントロピーの和が出来るだけ小さくなり, $H(p_c)$ と一致すれば cの各成分は独立であるといえる.以上より,独立成分分析の定式化の一つとして,

$$\min_{D} \|\boldsymbol{x} - D\boldsymbol{c}\|_{2}^{2} + \sum_{j=1}^{m} H(p_{c^{j}}), \qquad (26)$$

が得られる.ここで, $H(p_{c^j})$ は $c^j, j = 1, \ldots, m$ の密度関数で値が定まる汎関数であり, $c = D^{-1}x$ なので結局は D^{-1} の関数であることに注意する.原信号が優ガウス的な分布に従うときに,エントロピーは低い値を取る.その意味で,エントロピー最小化による独立成分分析は原信号にスパース性の仮定をおいて信号分離を行なっていると解釈できる.例えば観測信号がn 個ある時,係数ベクトルの第j成分がラプラス分布 $p(c^j) \propto \exp(-\lambda |c^j|)$ に従うとして,各成分の観測値を $c_i^j, i = 1, \ldots, n$ とすると, $p(c^j)$ に関する期待値を経験平均で置き換えることでエントロピー $H(p_{c^j})$ は

$$H(p_{c^{j}}) = -\frac{1}{n} \sum_{i=1}^{n} \log p_{j}(c_{i}^{j}) \propto \lambda \sum_{i=1}^{n} |c_{i}^{j}| = \lambda \|\boldsymbol{c}^{j}\|_{1} \quad (27)$$

のように推定できる.これは,Ψ(D,C)として行列の成分 毎の l₁ ノルムを用いることに相当する.以上より,問題 (26) は正則な行列 D に関する最適化問題

$$\min_{D} \|X - DC\|_{2}^{2} + \lambda \|D^{-1}X\|_{1}$$
(28)

になる.ただし,原信号と観測信号の間には $C = D^{-1}X$ なる関係があることに注意する.

5.4 非負行列因子分解

非負行列因子分解の場合,ポアソン分布に基づく式 (17) のような損失関数も用いられるが,ここでは簡単のため近 似誤差として正規分布を仮定した時の負の対数尤度で損失 関数 $L(D,C;X) = ||X - DC||_2^2$ を考える. D 及び C の各 成分には非負という制約があるため,これらの事前分布と しては非負の値を取る分布を設定する必要がある.ここで は一つの確率的解釈として,D, Cの各成分がそれぞれ独立 にパラメタ λ_1 と λ_2 の指数分布に従うとする.このとき,

$$p(D) = \prod_{i=1}^{d} \prod_{j=1}^{m} \lambda_1 \exp(-\lambda_1[D]_{ij}), \quad [D]_{ij} \ge 0,$$

$$p(C) = \prod_{j=1}^{m} \prod_{k=1}^{n} \lambda_2 \exp(-\lambda_2[C]_{jk}), \quad [C]_{jk} \ge 0$$

であり, D, C の MAP 推定は

$$\min_{[D]_{ij} \ge 0, [C]_{jk} \ge 0} \|X - DC\|_{2}^{2}
+ \lambda_{1} \sum_{i=1}^{d} \sum_{j=1}^{m} [D]_{ij} + \lambda_{2} \sum_{j=1}^{m} \sum_{k=1}^{n} [C]_{jk}
= \min_{[D]_{ij} \ge 0, [C]_{jk} \ge 0} \|X - DC\|_{2}^{2} + \lambda_{1} \|D\|_{1} + \lambda_{2} \|C\|_{1}
(29)$$

と等価である.さらに残差の分布及び [D]_{ij}, [C]_{jk}の従う 指数分布のパラメタにも事前分布を設定し, Gibbs サンプ リングにより Bayes 的に非負行列因子分解を行う手法が提 案されている [28].

6. 係数選択の手法

観測信号 x と辞書 D が与えられた時, x を Dc で近似 するような係数 c を求める問題を,(狭義の)スパースコー ディング問題と呼ぶ.ここでは最適化問題(3)を,再構成 誤差を一定の閾値以下に抑えた上で出来るだけ少ない数の 基底の線型結合で信号を近似する問題

$$\min_{\boldsymbol{c}} \|\boldsymbol{c}\|_0 \text{ subject to } \|\boldsymbol{x} - D\boldsymbol{c}\|_2 < \epsilon \qquad (30)$$

の形に書き直して考える.先に述べたように,スパース性の制約が l_0 ノルム制約の場合には (30) は組合せ最適化問題であり,NP 困難な問題である [10].この問題に対する解法として,貪欲法に基づく方法や, l_0 制約を l_1 制約で緩和した上で解く方法など,数多くのアルゴリズムが提案されている.本節では,信号処理及び画像処理の分野で広く用いられているスパースコーディングアルゴリズムとして,

(1) Orthogonal Matching Pursuit [29]

(2) Iterative Reweighted Least Squares [30] を紹介する.

Orthogonal Matching Pursuit (OMP)

OMP は,観測信号の近似に利用する係数の添字集合の中から「サポート」,すなわち非ゼロ係数の添字集合 S を見つけ出すアルゴリズムである.初めはサポートは空集合であるとして,観測信号 x を基底の線型結合で近似した時の残差を最小にするように新たな基底をサポート集合に一つ一つ追加していき,サポートに含まれる基底のみで信号を近似した時の残差が ϵ 以下になったら停止する.残差の低

減に寄与する基底を順次選択していく貪欲法であり,解の 最適性は保証されないが,多くの場合優れた近似を与える ことが知られている.OMPの手続きを Algorithm1 にま とめる.

Algorithm 1 OMP: (Orthogonal Matching Pursuit) input: 辞書行列 $D \in \mathbb{R}^{d \times m}$ 観測信号 $oldsymbol{x} \in \mathbb{R}^d$ 近似閾値 $\epsilon > 0$ initialize: 繰り返しのカウンタt=0初期係数ベクトル $m{c}^0 = m{0}$ 初期残差 $\boldsymbol{r}^0 = \boldsymbol{x} - D\boldsymbol{c}^0 = \boldsymbol{x}$ 初期サポート集合 $S = \text{supp}\{c^0\} = \emptyset$ while $\|\boldsymbol{r}^t\|_2 \geq \epsilon$ do 残差計算: $j \notin S$ に対して, 基底 d_i を追加した時の残差を計算: $\epsilon(j) = \min_{z_j} \|\boldsymbol{d}_j z_j - \boldsymbol{r}^t\|_2^2$ サポート更新: $j_0 = \arg\min_{j \notin S} \epsilon(j)$ をサポートに追加: $\mathcal{S} = \mathcal{S} \cup \{j_0\}$ 暫定解の計算: $c^{t+1} = \arg\min \|Dc - x\|_2^2$ subject to $\sup\{c\} \subset S$

残差更新: $r^{t+1} = x - Dc^{t+1}$. end while

Iterative Reweighted Least Square (IRLS)

IRLS は, $p \in [0,1]$ に対して l_p ノルムを重み付き l_2 ノルムで繰り返し近似することで, l_p ノルム正則化問題を近似的に解く手法である. t 回目のくり返しで得られている係数ベクトルを $c^t \in \mathbb{R}^m$ とする.これを用いて,重み行列を $W_t = \text{diag}(|c_1^t|^{1-p/2}, \cdots, |c_m^t|^{1-p/2})$ で定義する.対角成分が0の時は逆行列の対応する成分もゼロにすることにして W_t^{-1} を定義すると, $||W_t^{-1}c^t||_2^2 = ||c^t||_p^p$ となり, $||W_t^{-1}c||_2^2$ は c の l_p ノルムを模していることが分かる.そこで,問題

min $||W_t^{-1}\boldsymbol{c}||_2^2$ subject to $||\boldsymbol{x} - D\boldsymbol{c}||_2 = \boldsymbol{0}$ (31)

を考えて, ラグランジュの未定乗数法により最適解

$$\boldsymbol{c}^{t+1} = W_t^2 \boldsymbol{D}^\top (\boldsymbol{D} W_t^2 \boldsymbol{D}^\top)^\dagger \boldsymbol{x}$$
(32)

を得る.ただし, M^{\dagger} は行列Mの一般化逆行列を表す.これを,残差のノルム $\|r^t\|_2 = \|x - Dc^t\|$ が ϵ 以下になるまで繰り返す.IRLSによる係数最適化の手続きを Algorithm2にまとめる.

OMP を始めとするスパース表現のための係数選択手法 は広く用いられており,高速・高精度な計算方法や,理論 的に興味深い手法が現在も開発されている.例えば[31]で は,OMP を初期解として,指定した計算時間内で OMP による解の改善を木探索アルゴリズムに基づいて行う方法 IPSJ SIG Technical Report

end while

が提案されており,計算時間と近似精度のトレードオフを 利用者が調整できる様になっている.また,観測信号を基 底表現における係数の情報が欠損した不完全データをみな すことで,EMアルゴリズムの枠組みで係数推定を捉える 手法 [32] や,Bayes 推定の枠組みとして OMP を捉えた研 究 [33] も行われている.

7. 基底学習の手法

基底の線型結合による信号の表現において,辞書のデザ インは非常に重要である.離散コサイン変換や Fourier 変 換,wavelet [3],あるいは curvelet [34]のように予め決め られた手段に従って基底を用意しておく方法と,信号から 基底を学習する方法がある.一般に基底の学習には大きな 計算コストがかかるが,近似対象の信号と類似した特徴を 有する信号から学習した辞書を利用することで,スパース 性が高く誤差の少ない近似が得られることが期待できる.

初期の研究においては,基底の学習は勾配法によって行 われた [1].計算効率や収束性の観点からは改善の余地の大 きい方法であるが, Gabor wavelet 基底に似た局所的な基 底が得られている、その後,信号処理の観点から,Engan らにより k-means 法に基づく基底学習手法である Method of Optimal Directions [35] が提案された. さらに k-means 法の一般化として Aharon らにより K-SVD [36] と呼ばれ る基底学習アルゴリズムが提案された.K-SVD を含め多 くの基底学習アルゴリズムは,固定された基底のもとで係 数の最適化を行い,その係数を固定して基底を最適化する ことを繰り返す.このように基底と係数の交互最適化手法 を用いる K-SVD の基本原理を以下に示す.係数行列 C は 例えば OMP などを利用して求められているとする.基本 的な考え方は, l 番目の基底 d_l を更新する際, d_l を利用し ないで信号を近似した時の誤差を考え,この誤差を表現す る基底を新たな d1 とするというものである. 具体的な手 段としては , まず観測信号 $X = (oldsymbol{x}_1, \dots, oldsymbol{x}_n)$ の中で , 現在 の基底 d_l が表現に用いられているような信号の添字集合

$$\Omega_l = \{i \in \{1, \dots, n\} | [C]_{li} \neq 0\}$$

を求め, Ω_l に含まれる観測信号のみからなる観測信号行列,すなわち Ω_l に含まれる添字の列ベクトルからなるXの部分行列 X^{Ω_l} を構成する.更新対象の基底 d_l を利用しないで観測信号を近似した時の残差を

$$R_l = X^{\Omega_l} - \sum_{j \neq l} \boldsymbol{d}_j \boldsymbol{c}^j \tag{33}$$

とすると, $||X^{\Omega_l} - DC||_2^2 = ||R_l - d_l c^l||_2^2$ であり, 最適な 基底 d_l は残差をランク1で l_2 ノルムの意味で最良近似す れば得られることが分かる.この最良近似は,残差行列 R_l に特異値分解を施し,その第一左特異ベクトル u_1 を d_l と すれば良い.K-SVD は基底の学習アルゴリズムであるが, 基底学習の際に合わせて係数の修正も $c^j = \sigma_1 v_1$ で行う. これを全ての基底ベクトル $d_l, l = 1, \ldots, m$ について順次 行うことで辞書 Dを更新する.この基底の最適化と係数の 最適化を,予め定めた基準を満たすまで繰り返す.この手 続きを Algorithm3 にまとめる.この基底学習アルゴリズ ムと,OMP などの係数選択アルゴリズムとを繰り返し行 うことで,観測信号に適した辞書を学習することが出来る.

Algorithm 3 K-SVD Dictionary Learning Algorithm
input:
辞書 $D \in \mathbb{R}^{d imes m}$
観測信号 $X=(oldsymbol{x}_1,\ldots,oldsymbol{x}_n)$
係数行列 C
for $l = 1, \ldots, m$ do
添字集合 Ω_l を計算する
残差行列 $R_l = X^{\Omega_l} - \sum_{j eq l} oldsymbol{d}_j oldsymbol{c}^j$ を計算する
特異値分解 $R_l = U \Sigma V^ op$ を行い,
$oldsymbol{d}_l \leftarrow oldsymbol{u}_1, oldsymbol{c}_l \leftarrow \sigma_1 oldsymbol{v}_1$
で基底,係数を更新する
end for

K-SVD アルゴリズムは実問題に対して優れた性能を示 しており,その数学的なシンプルさ,計算効率の高さも相 まって,画像処理,音声信号処理におけるスパース表現の ための基底学習手法として広く用いられている.

8. 画像処理への応用

本稿の最後に,画像のスパース表現の応用例を幾つか挙 げる.基底展開による画像の表現の代表的な応用は画像圧 縮及びノイズ除去であるが,他にも以下のような応用が可 能である.

画像分離,画像修復[37,38]

2つの異なる画像が重なりあった画像 *x* が観測されると する.この時,2つの異なる画像がそれぞれ異なる辞書 D_a, D_b と係数 c_a, c_b で $D_a c_a, D_b c_b$ のように表現されてい るとすると,次の問題

 $\min_{\boldsymbol{c}_a, \boldsymbol{c}_b} \|\boldsymbol{c}_a\|_0 + \|\boldsymbol{c}_b\|_0 \text{ subject to } \|\boldsymbol{x} - D_a \boldsymbol{c}_a - D_b \boldsymbol{c}_b\|_2^2 \leq \epsilon$

を解くことで重なりあった画像の分離を行うことが出来る. この問題は Morphological Component Analysis (MCA) と 呼ばれ,特にテクスチャーと線素の分離,あるいは単純な 独立ノイズではなく,構造を持ったノイズにより劣化した 画像の修復などに応用されている.

超解像手法 [39,40]

まず,高解像度の学習用画像から辞書 D_H を学習してお き,その辞書に対してボケ,ダウンサンプリングによる擬 似的な劣化を施すことで低解像度画像を表現するための辞 書 D_L を作成する.観測した複数枚の低解像度画像を辞書 D_L を用いて表現し,非ゼロの係数を記録しておく.この 非ゼロの係数に対応する高解像度画像用の基底に係数をか けて足し合わせることで,高解像度の画像を得る.この超 解像手法の概念図と,この手法により得られた高解像度画 像の例を図2に示す.



スパース表現を利用したロバストな顔認識 [41,42] 顔画像による認証は古くからコンピュータビジョン及びセ キュリティの重要な研究課題であり,スパース表現を利用 した顔画像認証手法も多数提案されている.特に [41] で は,スパース表現を直接的に利用した顔認識手法を提案し ている.認証システムに登録した人物がH人いるとして, 各人物は複数の顔画像 $f_1^h, \ldots, f_{n_h}^h, h = 1, \ldots, H$ を登録し ているとする.これらの画像を基底として扱い,辞書

$$D = (\boldsymbol{f}_1^1, \dots, \boldsymbol{f}_{n_1}^1; \cdots; \boldsymbol{f}_1^H, \dots, \boldsymbol{f}_{n_H}^H)$$

を構成し,認証対象の画像をこの基底のスパースな線型結合で表現する.このときの結合係数は認証対象の人物の画像が辞書に含まれていれば,その人物の登録画像に対応する少数の係数のみが非ゼロの値を取ることが期待されるため,基底の結合係数から認証スコアを算出して認証を行うことが出来る.さらに[42]ではロバスト統計で用いられるような外れ値に対して頑健な損失関数を採用することで,より安定した認識を実現している.

9. 終わりに

本稿では,信号のスパース表現を行列の分解の形式で整 理し,その確率モデルとしての解釈を試みた.信号をス パース表現する際に重要な係数及び基底の学習アルゴリズ ムを紹介し,画像処理における幾つかの応用例を挙げた. 本稿では説明出来なかったが,スパースにサンプリングし た信号からの原信号の復元可能性 [43] や,スパース正則化 回帰問題におけるオラクル性 [44] など, 圧縮センシングや 統計学の分野ではスパース表現からの情報復元の可能性の 研究が盛んに行われている.機械学習においてここ10年で 大きく発展した Multiple Kernel Learning (MKL [45]) も, 多数のカーネル関数によって表現される特徴量から有用な ものをスパースな凸結合により選択する手法と捉えること が可能であり [46],画像処理にも応用されている.スパー ス表現,あるいは行列の低ランク・スパース分解は現在で も盛んに研究されている分野である.スパース表現の範疇 として取り扱える問題は非常に広く,画像処理の諸問題へ の有効なアプローチとしてもさらに発展が期待される.

参考文献

- B. A. Olshausen and D. J. Field: "Emergence of simplecell receptive field properties by learning a sparse code for natural images", Nature, **381**, pp. 607–609 (1996).
- [2] B. A. Olshausen and D. J. Field: "How Close Are We to Understanding V1?", Neural Computation, 17, pp. 1665–1699 (2005).
- [3] T. S. Lee: "Image representation using 2D Gabor wavelets", IEEE Trans. Pattern Anal. Mach. Intell., 18, 10, pp. 959–971 (1996).
- [4] G. Palm: "On associative memory", Biological Cybernetics, 36, pp. 19–31 (1980).
- [5] S. Amari: "Characteristics of sparsely encoded associative memory", Neural Networks, 2, 6, pp. 451–457 (1989).
- S. Akaho: "Capacity and error correction ability of sparsely encoded associative memory", ICANN'93, pp. 707–710 (1993).
- [7] S. Z. Li: "Markov random field modeling in image analysis", Springer-Verlag (2001).
- [8] K. Tanaka: "Statistical-mechanical approach to image processing", Journal of Physics A: Mathematical and General, 35, 37, pp. R81–R150 (2002).
- [9] M. E. Tipping and C. M. Bishop: "Bayesian image superresolution", NIPS15, pp. 1303–1310 (2003).
- [10] B. K. Natarajan: "Sparse approximate solutions to linear

- [11] M. Elad: "Sparse and Redundant Representations: From Theory to Applications in Signal and Image Processing", Springer-Verlag (2010).
- [12] D. Harville: "Matrix Algebra from a Statistician's Perspective", Springer-Verlag (2008).
- [13] I. Jolliffe: "Principal Component Analysis", Springer-Verlag (1986).
- [14] H. Zou, T. Hastie and R. Tibshirani: "Sparse principal component analysis", Journal of Computational and Graphical Statistics, 15, pp. 265–286 (2004).
- [15] A. Hyvärinen, J. Karhunen and E. Oja: "Independent Component Analysis", John Wiley Sons, Inc. (2001).
- [16] D. D. Lee and H. S. Seung: "Learning the parts of objects by non-negative matrix factorization", Nature, 401, 6755, pp. 788–791 (1999).
- [17] D. D. Lee and H. S. Seung: "Algorithms for non-negative matrix factorization", NIPS13, pp. 556–562 (2000).
- [18] P. O. Hoyer: "Non-negative matrix factorization with sparseness constraints", J. Mach. Learn. Res., 5, pp. 1457–1469 (2004).
- [19] A. Cichocki, H. Lee, Y.-D. Kim and S. Choi: "Nonnegative matrix factorization with alpha-divergence", Pattern Recognition Letters, 29, 9, pp. 1433–1440 (2008).
- [20] Y. Fujimoto and N. Murata: "Nonnegative matrix factorization via generalized product rule and its application for classification", LVA/ICA2012, pp. 263–271 (2012).
- [21] T. J. Mitchell and J. J. Beauchamp: "Bayesian Variable Selection in Linear Regression", Journal of the American Statistical Association, 83, 404, pp. 1023–1032 (1988).
- [22] S. Mohamed, K. Heller and Z. Ghahramani: "Bayesian and l1 approaches to sparse unsupervised learning", ICML2012, pp. 721–728 (2012).
- [23] M. K. Titsias and M. Lázaro-Gredilla: "Spike and slab variational inference for multi-task and multiple kernel learning", NIPS24, pp. 2339–2347 (2011).
- [24] M. E. Tipping and C. M. Bishop: "Probabilistic principal component analysis", Journal of the Royal Statistical Society, Series B, 61, pp. 611–622 (1999).
- [25] B. Everitt: "An Introduction to Latent Variable Models", Monographs on Statistics and Applied Probability, Chapman and Hall (1984).
- [26] Y. Guan and J. G. Dy: "Sparse probabilistic principal component analysis", J. Mach. Learn. Res. - Proceedings Track, 5, pp. 185–192 (2009).
- [27] T. M. Cover and J. A. Thomas: "Elements of information theory", John Wiley and Sons, Inc. (1991).
- [28] M. N. Schmidt, O. Winther and L. K. Hansen: "Bayesian non-negative matrix factorization", ICA2009, pp. 540– 547 (2009).
- [29] Y. C. Pati, R. Rezaiifar, Y. C. P. R. Rezaiifar and P. S. Krishnaprasad: "Orthogonal matching pursuit: Recursive function approximation with applications to wavelet decomposition", Asilomar1993, pp. 40–44 (1993).
- [30] R. Chartrand and W. Yin: "Iteratively reweighted algorithms for compressive sensing", ICASSP2008, pp. 3869– 3872 (2008).
- [31] R. Rei, J. P. Pedroso, H. Hino and N. Murata: "A tree search approach to sparse coding", LION6 (2012, to appear).
- [32] A. C. Gurbuz, M. Pilanci and O. Arikan: "Expectation maximization based matching pursuit", ICASSP2012, pp. 3313–3316 (2012).

- [33] A. Dremeau, C. Herzet and L. Daudet: "Structured Bayesian orthogonal matching pursuit", ICASSP2012, pp. 3625–3628 (2012).
- [34] E. Candès and D. Donoho: "Curvelets: A surprisingly effective nonadaptive representation for objects with edges", Curves and Surfaces (Ed. by L. L. Schumaker et al.), Vanderbilt University Press (1999).
- [35] K. Engan, S. O. Aase and J. Hakon Husoy: "Method of optimal directions for frame design", ICASSP1999, pp. 2443–2446 (1999).
- [36] M. Aharon, M. Elad and A. Bruckstein: "K-SVD: An algorithm for designing overcomplete dictionaries for sparse representation", IEEE Trans. Sig. Proc., 54, 11, pp. 4311–4322 (2006).
- [37] J. Bobin, J. L. Starck, J. M. Fadili, Y. Moudden and D. L. Donoho: "Morphological component analysis: An adaptive thresholding strategy", IEEE Trans. Img. Proc., 16, 11, pp. 2675–2681 (2007).
- [38] M. Elad, J. Starck, P. Querre and D. Donoho: "Simultaneous cartoon and texture image inpainting using morphological component analysis (MCA)", Applied and Computational Harmonic Analysis, **19**, 3, pp. 340–358 (2005).
- [39] Y. Ueda, T. Ohta, A. Iwase, M. Seki and N. Murata: "Super-resolution based on sparse coding", Proceedings of Computational Algebraic Statistics, Theories and Applications (2008).
- [40] J. Yang, J. Wright, T. S. Huang and Y. Ma: "Image super-resolution via sparse representation", IEEE Trans. Image Proc., 19, 11, pp. 2861–2873 (2010).
- [41] J. Wright, A. Y. Yang, A. Ganesh, S. S. Sastry and Y. Ma: "Robust face recognition via sparse representation", IEEE Trans. Pattern Anal. Mach. Intell., **31**, 2, pp. 210–227 (2009).
- [42] M. Yang, L. Zhang, J. Yang and D. Zhang: "Robust sparse coding for face recognition", CVPR2011, pp. 625– 632 (2011).
- [43] E. J. Candès and T. Tao: "Decoding by linear programming", IEEE Trans. Inf. Theo., 51, 12, pp. 4203–4215 (2005).
- [44] J. Fan and R. Li: "Variable selection via nonconcave penalized likelihood and its oracle properties", Journal of the American Statistical Association, 96, 456, pp. 1348– 1360 (2001).
- [45] A. Rakotomamonjy, F. R. Bach, S. Canu and Y. Grandvalet: "SimpleMKL", J. Mach. Learn. Res., 9, pp. 2491– 2521 (2008).
- [46] F. R. Bach: "Consistency of the group lasso and multiple kernel learning", J. Mach. Learn. Res., 9, pp. 1179–1225 (2008).