

独立成分解析とは

ノイズと独立成分解析

科学技術振興事業団

さきがけ研究 21

池田 思朗

花子: 先生, 日焼けしてますね. またスキーですか?

先生: おいおい, そんなに暇ではないぞ. 確かにこの冬もアメリカのデンバー, 北海道と出掛けたが, ちゃんと会議に出席していたんだ. そのついでに, ちょっとだけスキーも ...

次郎: なんだか怪しいけど, まあいいや.

一夫: 今日は次郎君から質問があると聞いているんだけど?

次郎: あ, はい. 実は前回の本田先生のお話が面白かったので, 僕も ICA を使って多次元の生体データを分離してみたのです. でも, なんだかうまくいかないんです.

先生: 次郎君は, ますます実践づいているな. ICA は生体信号や音声信号, その他にも様々な多次元信号解析への応用が注目されている. 実際に使ってみると新たな問題も見つかるだろう. 大変結構.

花子: それで, うまくいかないというのはどういうこと?

次郎: うん. 実は, MEG のデータを手に入れたので, この間話にあがっていた, A. Hyvärinen の FastICA, J-F. Cardoso の JADE, それから Bell and Sejnowski のアルゴリズムを web から拝借して, やってみたんだ.

花子: まあ, 色々なアルゴリズムを使ってみたのね.

次郎: うん. まずはやってみないと要領を得なくてさ. それに, 本当にどのアルゴリズムも信号を独立成分に分離するなら, どれも同じ結果になるはずだろう?

一夫: それで, 結果はどうだったんだい?

次郎: それで、アルゴリズム毎に違う結果になっちゃって、なんだか分らなくなっただ。

先生: なにやら興味深い話だね。もう少し詳しく実験の説明をしてくれ。

次郎: このデータは 126 チャンネルのデータ (10 次元のみを図 1 に示す) なんですけど、脳の計測データじゃなくて、模型に電極を埋めこんで、三角波を入力して測定したデータなんだ。だから、実際のデータはなんだか分っているんだ。データは時間方向には数百点ほどのデータだから、このまま ICA を掛ければ、126 個の独立成分に分離できるんだけど、web から拝借してきたのは、どれも次元が多くなるとアルゴリズムがなかなか収束しなくて、うまくいかないんだ。

花子: でも、さっき結果が全部違うっていったじゃない? 結果は出たんでしょう? 結果を計算するところまではどうやったの?

次郎: まあ、待ってよ。まだ、話に続きがあるんだよ。

アルゴリズムが収束しないのは次元が高いからだと思ったんだ。だから、まず次元を減らすために PCA で処理したんだ。共分散行列の固有値の大きい方からいくつか次元を決めて固有ベクトルを選んで、その固有ベクトルの空間に射影することで、次元を落したんだ。試しに 10 次元に落したんだけど、そのあと ICA で処理すると、アルゴリズム毎に違う結果になってしまうんだ (図 2~4)。

先生: なるほど大体状況は分った。一夫君、どんな可能性が考えられるかな?

一夫: はい。その前に聞きたいんだけど、何かデータについて、聞いたことはなかったかい?

次郎: MEG データは各センサに独立にノイズがあるから、100 回程度の加算平均が必要だと言う話です。

一夫: 多分そのノイズが影響しているんだろうな。MEG などの生体計測データにはどのようなノイズがあると思う?

花子: 今の話だと各センサーに独立に加わっているノイズがあるわね。それから、脳の計測データは微弱な信号を観測するから、環境から全てのセンサーに影響するノイズがあるでしょうね。

一夫: そうだね。多分モデルとしては、

$$x = As \tag{1}$$

ではなくて、

$$\boldsymbol{x} = \boldsymbol{A}\boldsymbol{s} + \boldsymbol{\epsilon} \quad (2)$$

の方が生体計測データには近いのだろうな。 \boldsymbol{s} は興味のある信号と、全てのセンサーに加わっているノイズ、そして $\boldsymbol{\epsilon}$ は各センサーに独立に加わるノイズというわけだ。

花子: わかったわ。私達は信号を、独立成分とノイズと区別しているけど、アルゴリズムはそれらを必ずしも区別しないから、

$$\boldsymbol{x} = \boldsymbol{A}\boldsymbol{s} + \boldsymbol{\epsilon} = \boldsymbol{A}'\boldsymbol{s}' = (\boldsymbol{A}, \boldsymbol{I}) \begin{pmatrix} \boldsymbol{s} \\ \boldsymbol{\epsilon} \end{pmatrix} \quad (3)$$

のように書き直せば、信号の次元を m とすれば、 \boldsymbol{A}' は $n \times (m+n)$ の行列になって、このままでは信号源の数の方が観測よりも多くなってしまふのね。

次郎: でも、信号源の数の方が観測の数よりも多いと、線型な行列の解は存在しないよ。

花子: それに、 $\boldsymbol{\epsilon}$ で表されるノイズには興味はないわね。推定したいのは \boldsymbol{A} と \boldsymbol{s} だけだもの。

先生: 確かにノイズの問題は実際のデータを ICA で処理するときには問題になる。ここでまず 1 回考えてみるのも良いだろう。MEG のデータの解析については面白い結果が報告されているが、結果についてはまた別の機会に皆で調べよう。

では、次郎君、なぜアルゴリズムがうまくいかなかったのか、説明してくれ。

次郎: はい。このモデルが正しいとすれば、まず PCA の部分がうまく働かなかったんだと思います。次元を減らす際に用いた PCA もうまく働かないし、ICA の pre-whitening または sphering と呼ばれる前処理でも PCA を用いているわけで、この部分も正しく働かないんでしょう。

花子: うまく働かないってどういうこと? もう少し具体的に話してよ。

次郎: うん。主成分分析は信号のパワー、つまり 2 次の統計量のみを用いるんだから、共分散行列のみを考えれば十分だろう。だから、この場合、ノイズの含まれた観測信号の共分散行列を考えればすぐに分るよ。さっきの式で \boldsymbol{x} の共分散を \boldsymbol{C} とするよね。 $E[\cdot]$ を平均として、 $\boldsymbol{C} = E[\boldsymbol{x}\boldsymbol{x}^T]$ だ。仮に信号の共分散行列が m 次元の単位行列で \boldsymbol{I}_m だということも分かっているとしよう。すると、

$$\boldsymbol{C} = \boldsymbol{A}\boldsymbol{A}^T + \boldsymbol{\Sigma}$$

だろう？もし、なんらかの方法で混合行列 A に対する分離のための行列 W が求まったとしよう。

花子: どうやるかは別にして、まあ良いわ。

次郎: 今、信号のパワーは分っているから、大きさの任意性はなく、2つの行列の積 WA は、順序の入れ換えの任意性しかもたないわけだ。すると、 $y = Wx$ の共分散行列は、

$$E[yy^T] = E[Wxx^TW^T] = I_m + W\Sigma W^T$$

となってしまうと、分離された信号は必ずしも無相関にはならないんだ。答えが無相関じゃないのだから、逆に言えば観測信号を PCA で無相関化しても、前処理としての役目を果たさないということでしょう。

花子: 本当は信号部分のみを無相関にしたいのに、ノイズも含めて無相関化してしまうということなのね。

先生: 問題点は分ったとして、今、興味があるのは A と s だけで ϵ には興味がないんだから、それらを推定することを考えようじゃないか。

次郎: ICA での仮定は信号が独立で平均が 0 ということだったけど、この場合 s と ϵ のそれぞれの成分が互いに独立で平均が 0 だというだけだとうまく行かないでしょう。なにか他に仮定をしないと。

先生: そうだな。それでは、ノイズ ϵ が多次元の正規分布に従うとしようじゃないか。ついでにノイズの共分散行列が Σ と分っているとしたらどうだね？

次郎: それならば簡単ですよ。正規分布で共分散行列も分っているとしたら、ノイズに関しては全て分っているわけだから、ノイズを観測から差し引けば普通の ICA の問題になります。

花子: ノイズについて全部分ったっていったて、確率分布が分っただけで、 ϵ が分ったわけではないわ。

次郎: そうだけど、例えば x の共分散行列を C とするよね。 $C - \Sigma$ を計算すれば、 $E[Ass^T A^T] = C - \Sigma$ だから、ノイズを除いた共分散行列が推定できる。 s の次元 m が観測の次元 n より少ないとすれば、 $C - \Sigma$ の次元は m だから、これから信号の次元も推定できるし、他のキュムラントのような高次統計量には ϵ は何も影響しないはずだから、JADE などの高次統計量に基づく ICA のアルゴリズムは影響を受けないし、何もかも分ったも同然じゃない？

先生: C は推定しなければならないから, そう理想的に話は進まないかもしれないがな. では, ノイズが正規分布だとして, 共分散行列が分らないとしたら?

次郎: 共分散行列が分らなかったら, それを推定すれば良いのかな.

先生: そうだね. では, そのような手法について誰か知っているものは居ないのかな.

一夫: H. Attias, が提案した Independent Factor Analysis (以下, IFA) という手法がありますね.

先生: さすがは一夫君. それではざっと説明してくれるかな.

一夫: IFA では,

$$\boldsymbol{x} = A\boldsymbol{s} + \boldsymbol{\epsilon} \quad (4)$$

というモデルで A を推定する問題を考えるんだ. $\boldsymbol{\epsilon}$ は正規分布に従っているとしている. ただし, 論文では $\boldsymbol{\epsilon}$ は互いに相関のあるモデルを使っているから, 先生の言っていたのとは違って, 共分散行列は対角行列ではない. つまり, ノイズの分布は次のようになっている. 共分散行列を Λ として, ノイズの分布密度関数は,

$$p(\boldsymbol{\epsilon}; \boldsymbol{\theta}^\epsilon) = \mathcal{G}(\boldsymbol{\epsilon}; \Lambda) = \frac{1}{(2\pi|\Lambda|)^{n/2}} \exp\left(-\frac{1}{2}\boldsymbol{\epsilon}^T \Lambda^{-1} \boldsymbol{\epsilon}\right). \quad (5)$$

ただし, $\boldsymbol{\theta}^\epsilon = \{A_{ij}\} (i \leq j)$.

ここから IFA の独自のやり方なんだけど, さらに \boldsymbol{s} にもモデルを仮定したんだ. Attias は \boldsymbol{s} の各成分が混合正規分布に従うとしてモデルを作ったんだ. 各信号の分布は,

$$p(s_i; \boldsymbol{\theta}^i) = \sum_{j=1}^{L_i} \omega_{ij} \frac{1}{\sqrt{2\pi\sigma_{ij}^2}} \exp\left(-\frac{(s_i - \mu_{ij})^2}{2\sigma_{ij}^2}\right), \quad i = 1, \dots, m.$$

としたんだ. ここでも, $\boldsymbol{\theta}^i = (\omega_{i1}, \dots, \omega_{iL_i}, \mu_{i1}, \dots, \mu_{iL_i}, \sigma_{i1}^2, \dots, \sigma_{iL_i}^2)^T$ で, $\sum_j \omega_{ij} = 1, \omega_{ij} \geq 0$ が条件となる.

すると, $p(\boldsymbol{x}; \boldsymbol{\theta})$ が $\boldsymbol{\theta} = \{\boldsymbol{\theta}^1, \dots, \boldsymbol{\theta}^m, \boldsymbol{\theta}^\epsilon\}$ を使って

$$\begin{aligned} p(\boldsymbol{x}; \boldsymbol{\theta}) &= \int p(\boldsymbol{x}|\boldsymbol{s}; \boldsymbol{\theta})p(\boldsymbol{s}; \boldsymbol{\theta})d\boldsymbol{s} \\ &= \int \mathcal{G}(\boldsymbol{x} - \boldsymbol{s}; \Lambda) \prod_{i=1}^m p(s_i; \boldsymbol{\theta}^i) d\boldsymbol{s} \end{aligned} \quad (6)$$

と書けるから, これを使ってデータから最尤法で解けば全て求められるんだ.

次郎君: なるほど. 1つ1つの信号源の確率分布を L_i 個の正規分布からなる混合正規分布で近似したんだ. でもこれでうまくいくのかなあ. 少なくとも, 積分記号の中の \prod をなんとかしないと良く分らないよ. $p(s_i; \theta^i)$ は混合正規分布なんだから, 正規分布の重み付きの足し算になっているはずだ. それならば, 足し算を展開して, 個々の項を丁寧に拾い出せば, もう少し見易くなるだろう. $r = (r_1, \dots, r_m)^T$ を m 次元のベクトルとして, 各成分 r_i はそれぞれ $1, \dots, L_i$ の値を取るものとしよう. 可能な全ての r に対して加え合わせれば \prod を展開できる.

$$\omega_r = \prod_{i=1}^m \omega_{ir_i}, \boldsymbol{\mu}_r = (\mu_{1r_1}, \dots, \mu_{mr_m})^T, V_r = \text{diag}(\sigma_{1r_1}^2, \dots, \sigma_{mr_m}^2) \quad (7)$$

と定義すれば,

$$p(s; \theta^i) = \prod_{i=1}^m p(s_i; \theta^i) = \sum_r \omega_r \mathcal{G}(s - \boldsymbol{\mu}_r; V_r)$$

となる. 信号とノイズは独立だから,

$$\begin{aligned} p(\mathbf{x}; \boldsymbol{\theta}) &= \int \mathcal{G}(\mathbf{x} - \mathbf{s}; \Lambda) \sum_r \omega_r \mathcal{G}(s - \boldsymbol{\mu}_r; V_r) ds \\ &= \sum_r \omega_r \int \mathcal{G}(\mathbf{x} - \mathbf{s}; \Lambda) \mathcal{G}(s - \boldsymbol{\mu}_r; V_r) ds \\ &= \sum_r \omega_r \mathcal{G}(\mathbf{x} - A\boldsymbol{\mu}_r; AV_r A^T + \Lambda) \end{aligned} \quad (8)$$

となって, 全体が一つの大きな混合正規分布となる. r は $\prod_{i=1}^m L_i$ の取り方があるから, これだけの数の正規分布の足し算になるんだね.

あとは最尤推定ならば, 観測データを x_1, \dots, x_N として,

$$L = \frac{1}{N} \sum_{l=1}^N \log p(x_l; \boldsymbol{\theta})$$

を最大にするようにパラメータを決めれば良いんじゃないかな?

先生: そこまではそれで良いだろうが, では, 具体的にはどうやって解くん
だい?

次郎: 陽には解けそうにないから, やっぱ最急降下法になるのでしょうかね.
微分をとってグラジェントの方向にパラメータを動かしていけば解けると思
うんですが. 全てのパラメータがそれぞれ関係しているから, 多少面倒そう
だな.

一夫: それでも良いけど, Attias は EM(Expectation-Maximization) アルゴ
リズムによる推定法を導いている. EM アルゴリズムは Dempster らの 1977

の論文が元となっているが、直接観測できない確率変数を含むモデルの最尤推定のためのアルゴリズムなんだ。EM アルゴリズムの特徴は、比較的簡単な E-step と M-step を繰り返すことで最終的に最尤推定のためのパラメータが推定できるというものなんだ。

花子: EM アルゴリズムなら知っています。この場合、直接観測できない確率変数は r ということになるのかしら。それならば、現在のパラメータを θ_t として、次の時点のパラメータ θ_{t+1} は、次の E-step と M-step を通じて得られます。

E-step $Q(\theta, \theta_t)$ を次のように定義します。

$$Q(\theta, \theta_t) = \frac{1}{N} \sum_{l=1}^N \sum_r p(x_l | r; \theta_t) \log p(x_l, r; \theta)$$

ただし、

$$p(x | r; \theta) = \mathcal{G}(x - A\mu_r; AV_r A^T + \Lambda)$$

$$p(x, r; \theta) = \omega_r \mathcal{G}(x - A\mu_r; AV_r A^T + \Lambda)$$

です。

M-step は、

$$\theta_{t+1} = \operatorname{argmax}_{\theta} Q(\theta, \theta_t)$$

として求めます。これを繰り返せば、最終的には尤度の極大点が求まるということです。

先生: その通りだね。最終的な更新則を得るには、パラメータ間の制約を考慮に入れる必要があるから、もう少し計算が必要だね。Scaling の不定性もなんらかの方法で決めておかないと、収束しなくなってしまうしね。そのあたりは自分でできるだろうし、興味があれば書き下してみれば良いだろう。最後に文献を挙げておこう。

次郎: でもなんだか大変そうだな。パラメータの数がすごく増えているし、 $p(x; \theta)$ 全体の混合正規分布の正規分布数が、それぞれの信号源の分布で用いている混合正規分布の数の組み合わせで増えているから、平均操作も手間取りそうだ。これでは、MEG なんかの多次元の実際のデータではうまく動くのかよく分らないな。

花子: ICA の面白いところは、信号源の分布が良く分からなくても分離行列 W は推定できることだったと思うんだけど、これでは全てをパラメータで表現

しているから，ICA で陽には扱わなかった信号源の分布の部分まで扱わなくてはいけなくなってしまうわ．

一夫: ICA のセミパラメトリックモデルという側面は確かに無いね．でも，それはそれとして，面白い手法だと思うよ．もちろん，具体的に MEG などの高次元のデータに用いる場合には，混合分布の個数や次元 m の推定など，解決しなければいけない問題も残るけどね．

次郎: でも，もっと簡単にノイズを推定し，ノイズの影響を無くす手法は他にないの？

先生: もっと簡単に，か．それを考えるのなら，もう一度ノイズを含む観測のモデルから，問題点を考えてみよう．

$$x = As + \epsilon \quad (9)$$

ICA を実際に解く場合には，PCA で無相関化し，その後で ICA によって回転方向を決めるという手法を良くとる．その PCA の部分がうまく働かなかったんだから，PCA の部分だけを改善してみるのも 1 つの考え方だろう．

PCA はパワーつまり 2 次の統計量だけに注目していたんだから，この式の 2 次の統計量を考えてみよう．最初に言ったように信号 s の共分散行列が m 次元の単位行列 I_m ， ϵ は正規分布に従うとして，共分散行列が n 次元の対角行列 Σ だとしてみるんだ．

次郎: そうか． x の 2 次の統計量 $C = E[xx^T]$ から A と Σ を推定すれば良いんだ．モデルから得られる共分散行列は $AA^T + \Sigma$ だから，これらを一致させるように A と Σ を推定すれば良いんだね．

一夫: これは因子分析のモデルですね．因子分析は古くから研究のされてきた分野で，様々な手法があるけれど，次郎君の言うように，共分散行列のみから A と Σ と m を推定することが目的なんだ．因子分析の特徴を見ていこう．

まず，因子分析モデルは主成分分析と同様に回転の不定性を持っている．つまり， s を R によって $s' = Rs$ と直交変換したとすると， $As + \epsilon$ と $As' + \epsilon$ の 2 次の統計量は共に $AA^T + \Sigma$ だから， R をどのように取っても因子分析ではそれを推定できないことになる．つまり，回転に関しては一意に決まらない．

それから，自由度を考えると， C の自由度は $n(n+1)/2$ だね．一方， A と Σ のパラメータの数は $(m+1)n$ ．ただし，直交行列分の任意性を差し引かないと行けないので $m(m-1)/2$ を引くと，この自由パラメータが C の自由度よりも少ないということから，次の条件が導ける．

すなわち,

$$\frac{n(n+1)}{2} \geq (m+1)n - \frac{m(m-1)}{2}$$

$m < n$ という条件を加えてこれを解くと

$$m \leq \frac{1}{2} \{2n+1 - \sqrt{8n+1}\}.$$

となる. これは解が存在するための必要条件で Lederman によって 1937 に示されている.

また, 良く知られた A, Σ の推定法としては, 最小二乗法を用いて,

$$L(A, \Sigma) = \text{tr}(C - (AA^T + \Sigma))^2$$

を最小化する手法や, 信号に正規性を仮定し, 尤度を

$$L(A, \Sigma) = -\frac{1}{2} \{ \text{tr}(C(\Sigma + AA^T)^{-1}) + \log(\det(\Sigma + AA^T)) + n \log 2\pi \}$$

と定義し, 最尤推定する手法などがある.

次郎: では, 信号源の次数の推定はどうやるの.

一夫: 観測からは真の共分散は分からないから, やはり次数の推定も必要となるね. これにも様々な手法が研究されている. C の固有値によって推定するものや, 最尤推定したパラメータからモデル選択の手法を用いて AIC (Akaike Information Criteria) や MDL (Minimum Description Length) を用いて決定する手法などがある.

次郎: なるほど, いずれにせよ, 因子分析の手法を用いれば, ノイズの影響を考えた上で, ICA の前処理が可能になりそうだな. A の一般逆行列を使って信号を前処理すれば良いんだ. 一般逆行列を Q としよう. $AQA = A$ だから, $y' = Qx$ とすると, 正しく推定できたなら,

$$E[y'y'^T] = QAA^TQ^T + Q\Sigma Q^T = I_m + Q\Sigma Q^T$$

となるはずで, 信号部分は正しく前処理できるんだ.

一夫: もう一つ気をつけるのなら, 一般逆行列には不定性がある. 例えば $Q = (A^T \Sigma^{-1} A)^{-1} A^T \Sigma^{-1}$ と選べば回転の部分以外は y は s の最良線型不偏推定量となる.

次郎: そうか. それで, この Q で前処理をした信号を使って, 高次統計量に基づく JADE のような ICA の処理をすれば, ノイズの影響を受けずに ICA の処理が出来そうだ. 早速さっきのデータでやってみよう. 今回は最尤推定

を使ってパラメータを推定して、MDL で次元の推定をしてみよう。この場合は3次元と推定されたから、一般逆行列で処理をして、さらに JADE で処理をすれば、図5のようになりましたよ。

先生: 信号も雑音成分と、目的の信号と、なんとかそれらしい結果が得られたね。また別の機会にノイズについて取り上げるかもしれないが、今日はこのくらいにしておこう。

謝辞

MEG のデータは島津製作所から頂きました。ここに感謝します。

参考文献

- [1] Attias, H. (1999). Independent factor analysis. *Neural Computation*, 11(4), 803–851.

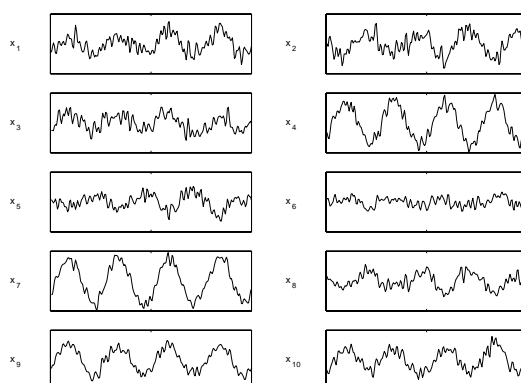


図 1:

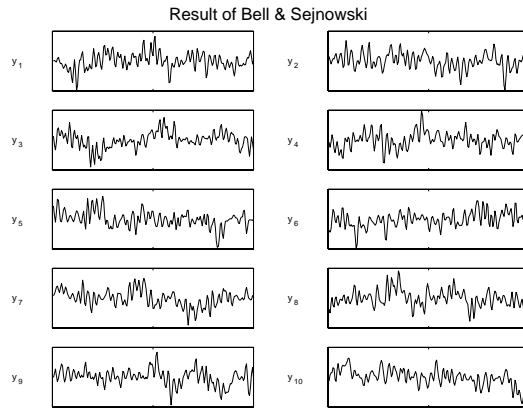


图 2:

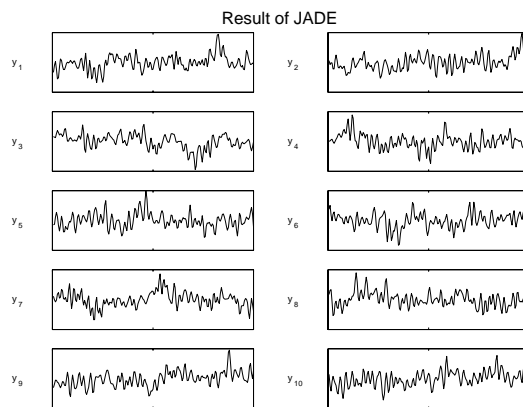
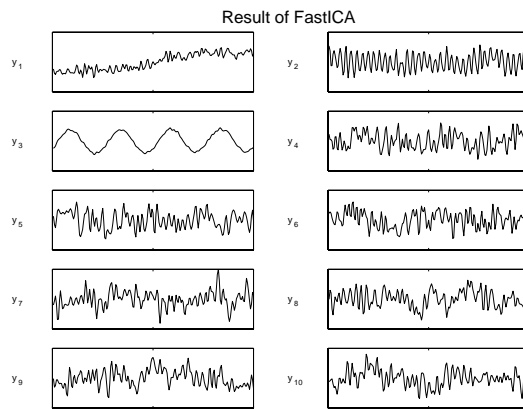
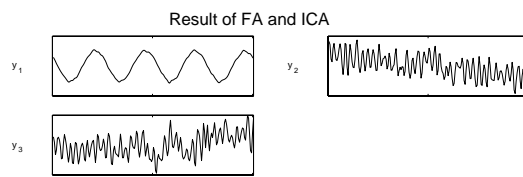


图 3:



☒ 4:



☒ 5: