# On the identifiability of ICA under multiple Gaussian sources

**started December 7, 2009; rewritten January 2, 2013**

**Gustavo Lacerda**

## Abstract

It is well-known that the Independent Component Analysis (ICA) model is not identifiable when more than one source is Gaussian (Comon 1993). In this note, we show that the number of continuous degrees of freedom in the "identifiable quotient space" of ICA is $\max(0, G - 1)$, where $G$ is the number of Gaussians in the sources.

However, since the source distributions are unknown, $G$ is unknown and needs to be estimated. We frame this as a problem of "hyper-equatorial regression" (i.e. find the hyper-great-circle that best fits the hyperspherical band), which is solved by doing PCA and discarding components beyond the knee of the eigenspectrum.

We conclude by suggesting a simple modification of the FastICA algorithm that only returns the identifiable components, by exiting as soon as no sufficiently non-Gaussian component can be found. This is more efficient than the first method because it avoids the bootstrap and clustering steps.

## 1 Independent Component Analysis (ICA)

The Noise-Free Independent Component Analysis [1, 2] is the following statistical model:

$\mathbf{X} = A\mathbf{S}$ (or in Einstein notation $\mathbf{X}_{ik} = a_i^j \mathbf{S}_{jk}$)

$\mathbf{S}_{jk} \sim Q_j$, i.i.d. for $k = 1, ..., K$

In other words, for each source $j$, we get a sample of $K$ i.i.d. observations from distribution $Q_j$. We may call the sample $\mathbf{S}_j$ the "realized source". $\mathbf{X}$ is observed ("mixtures"), and the $\mathbf{S}_j$ are assumed to be independent random variables with zero mean and unit variance ("sources").

The parameter of interest is the mixing matrix $A$. The distributions $Q_j$ are considered nuisance parameters.

It is well-known that $A$ can be identified up to column-sign (because simultaneously flipping $Q_j$ and negating the $j$th column of $A$ doesn't change the distribution of $X$) and row-permutations (because the $Q_j$ don't have any specified order), as long as $A$ is non-singular and *at most* one of the distributions in $S_i$ is Gaussian [1].

This note is concerned with another type of unidentifiability: the one that results from more than one source being Gaussian. In this case, rather than a discrete set of mixing matrices, there is now a continuous subspace of mixing matrices that result in the same observed joint distribution.

The equivalence relation yields a partition of the parameter space, i.e. a set of observationally-equivalent subspaces; this set of subspaces is called the "identifiable quotient space" or the "identifiable parameter".

If we have 2 Gaussian sources out of a large number of sources, ICA is far from useless: rather, all but these 2 components can be identified. Each member of the identifiable quotient space only has one degree of freedom (i.e. we have only lost 1 degree of freedom).

Below, we illustrate the behavior of the ICA statistic on an example where $I = J = 3$, that is, 3 sources and 3 mixtures. We resample about 500 times and plot the directions of the resulting components (columns of $A$) on the unit sphere:



(a) $G = 1$: 6 significant clusters are observed  (b) $G = 2$: 2 compact clusters and 1 ring-shaped cluster  (c) $G = 3$: no significant clusters
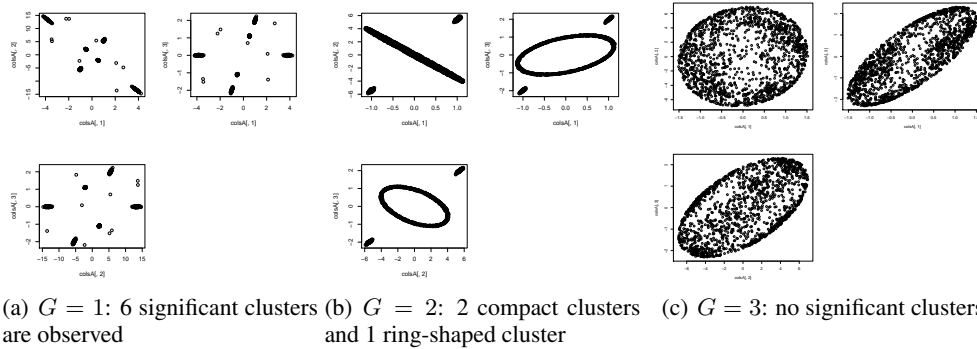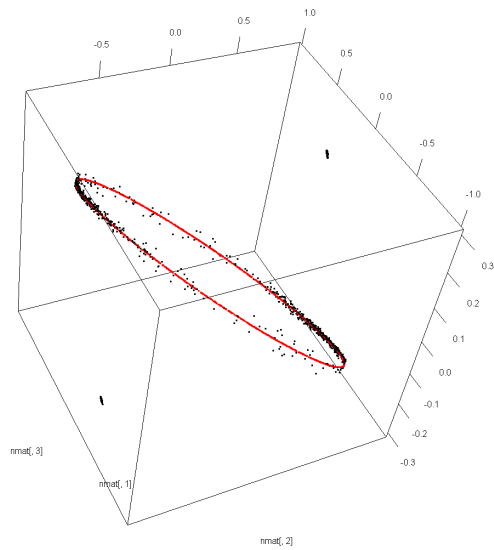
Figure 1: As we can see, the rank is $k - 1$

As we can see, in Fig (a), all 3 components have a clear direction, through their orientation can flip. This results in $2I = 6$ clusters. Fig (b) has one component like that, and two others that define a plane in which any pair of basis vectors fit the data equally well. In Figure (c) all components are Gaussian, so all vector triples are roughly equivalent.

What this suggests about ICA is that in general, rather than estimate single components, our goal should be to estimate subspaces.

The observational equivalence class that contains the truth is the set of mixing matrices that can be obtained by applying an invertible linear transformation of the columns of $A$ corresponding to the Gaussian components while maintaining the remaining columns fixed.

(a) $G = 2$

Figure 2: The tight clusters shown at the top right and bottom left of the figure indicate that the first component is stable. The red line traces the great circle that represents the best-fitting plane for the points around the band (by squared Euclidean distance), found using PCA.

Having run a clustering procedure to eliminate the tight clusters corresponding to the non-Gaussian components, we use PCA to find the k-plane around which the points on the ring lie. This gives us the Gaussian subspace.

The output of this modified ICA is now a pair: (a) a set of non-Gaussian directions (b) a k-plane representing the subspace defined by the Gaussian sources.

## References

[1] P. Comon (1994) - Independent component analysis - a new concept? *Signal Processing*, **36**:287-314.

[2] A. Hyvärinen, J. Karhunen, E. Oja (2001) - *Independent Component Analysis.* Wiley Interscience.