

A Tutorial on Independent Component Analysis

Jonathon Shlens*

Google Research
Mountain View, CA 94043

(Dated: April 14, 2014; Version 1.0)

Independent component analysis (ICA) has become a standard data analysis technique applied to an array of problems in signal processing and machine learning. This tutorial provides an introduction to ICA based on linear algebra formulating an intuition for ICA from first principles. The goal of this tutorial is to provide a solid foundation on this advanced topic so that one might learn the motivation behind ICA, learn why and when to apply this technique and in the process gain an introduction to this exciting field of active research.

I. INTRODUCTION

Measurements often do not reflect the very thing intended to be measured. Measurements are corrupted by random noise – but that is only one piece of the story. Often, measurements can not be made in isolation, but reflect the combination of many distinct sources. For instance, try to record a person’s voice on a city street. The faint crackle of the tape can be heard in the recording but so are the sounds of cars, other pedestrians, footsteps, etc. Sometimes the main obstacle preventing a clean measurement is not just noise in the traditional sense (e.g. faint crackle) but independent signals arising from distinct, identifiable sources (e.g. cars, footsteps).

The distinction is subtle. We could view a measurement as an estimate of a single source corrupted by some random fluctuations (e.g. additive white noise). Instead, we assert that a measurement can be a combination of many distinct sources – each different from random noise. The broad topic of separating mixed sources has a name - *blind source separation* (BSS). As of today’s writing, solving an arbitrary BSS problem is often intractable. However, a small subset of these types of problem have been solved only as recently as the last two decades – this is the provenance of *independent component analysis* (ICA).

Solving blind source separation using ICA has two related interpretations – filtering and dimensional reduction. If each source can be identified, a practitioner might choose to selectively delete or retain a single source (e.g. a person’s voice, above). This is a filtering operation in the sense that some aspect of the data is selectively removed or retained. A filtering operation is equivalent to projecting out some aspect (or dimension) of the data – in other words a prescription for dimensional reduction. Filtering data based on ICA has found many applications including the analysis of photographic images, medical signals (e.g. EEG, MEG, MRI, etc.), biological assays (e.g. micro-arrays, gene chips, etc.) and most notably audio signal processing.

ICA can be applied to data in a naive manner treating the technique as a sophisticated black box that essentially performs “magic”. While empowering, deploying a technique in this manner is fraught with peril. For instance, how does one judge the success of ICA? When will ICA fail? When are other methods more appropriate? I believe that understanding these questions and the method itself are necessary for appreciating when and how to apply ICA. It is for these reasons that I write this tutorial.

This tutorial is not a scholarly paper. Nor is it thorough. The goal of this paper is simply to educate. That said, the ideas in this tutorial are sophisticated. I presume that the reader is comfortable with linear algebra, basic probability and statistics as well as the topic of principal component analysis (PCA).¹ This paper does not shy away from informal explanations but also stresses the mathematics when they shed insight on to the problem. As always, please feel free to email me with any comments, concerns or corrections.²

II. MEASUREMENTS IN THE REAL WORLD

While the mathematics underlying the problems of blind source separation (BSS) are formidable, the ideas underlying BSS are quite intuitive. This section focuses on building this intuition by asking how signals interact in the real world when we make a measurement in which multiple sources contribute in addition to noise.

Many experiments have domain-specific knowledge behind how a measurement is performed. (How does a gene chip work?) Instead, this tutorial describes how measurements are made in domains in which we are all experts – our own senses.

¹ A good introduction to PCA written in the same style is *A Tutorial on Principal Component Analysis* by yours truly. This tutorial provides a full introduction and discussion of PCA with concrete examples for building intuition.

² I wish to offer a special thanks to E. Simoncelli for many fruitful discussions and insights and U. Rajashekar for thoughtful discussion, editing and feedback.

*Electronic address: jonathon.shlens@gmail.com

We discuss what a BSS problem entails in the visual and auditory domains and in the process build some intuition for what it takes to solve such a problem.

No better example exists of a BSS problem than sound. A fundamental observation from physics is that sound adds linearly. Whatever is recorded in a microphone or registered in our ear drum is the sum of the pressure waves from emanating from multiple sources. Consider an everyday problem. Amongst a background of music and conversations, we must discern a single person’s voice. The *cocktail party problem* is a classic problem in auditory signal processing. The goal of the cocktail party problem is to discern the sound associated with a single object even though all of the sounds in the environment are superimposed on one another (Figure 1). This is a challenging and highly applicable problem that only in recent years have reasonable solutions been devised based on ideas from ICA.

Image processing provides challenging BSS problems as well. Consider the problem of removing blur from an image due to camera motion (Figure 2). A photographer tries to take a photo but their camera is not steady when the aperture is open. Each pixel in the sensor array records the sum of all light within an integration period from the intended image along the camera motion trajectory. Each point in the recorded image can be viewed as the temporally weighted sum of light from the original image along the camera motion trajectory. Thus, each recorded pixel is the linear mixture of multiple image pixels from the original image. De-blurring an image requires recovering the original image as well as the underlying camera motion trajectory from the blurry image (Fergus *et al.*, 2006). Both the cocktail party and de-blurring problems are ill-posed and additional information must be employed in order to recover a solution.

These examples from vision and audition highlight the variety of signal interactions in the real world. No doubt more complex interactions exist in other domains or in measurements from specific instruments. Recovering individual sources from combined measurements is the purview of BSS and problems where the interactions are arbitrarily complex, are generally not possible to solve.

That said, progress has been made when the interactions between signals are simple – in particular, *linear* interactions, as in both of these examples. When the combination of two signals results in the superposition of signals, we term this problem a linear mixture problem. *The goal of ICA is to solve BSS problems which arise from a linear mixture.* In fact, the prototypical problem addressed by ICA is the cocktail party problem from Figure 1. Before diving into the mathematics of ICA it is important to visualize what linear mixed data “looks like” to build an intuition for a solution.

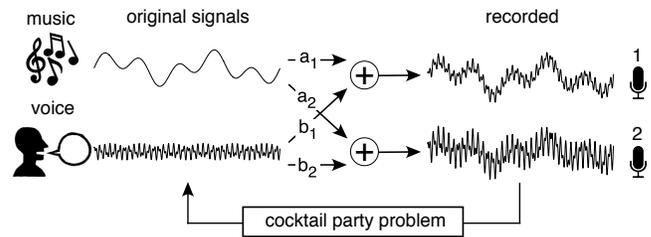


FIG. 1 Example of the cocktail party problem. Two sounds s_1, s_2 are generated by music and a voice and recorded simultaneously in two microphones. Sound adds linearly. Two microphones record a unique linear summation of the two sounds. The linear weights for each microphone (a_1, b_1 and a_2, b_2) reflect the proximity of each speaker to the respective microphones. The goal of the cocktail party problem is to recover the original sources (i.e. music and voice) solely using the microphone recordings (Bregman, 1994).

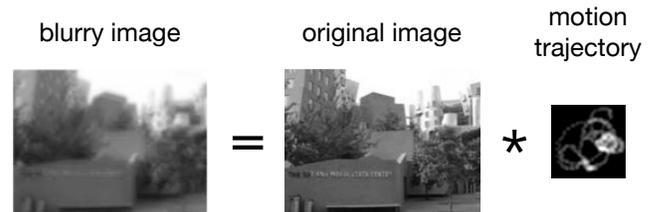


FIG. 2 Example of removing blur from an image due to camera motion. A blurry image (left panel) recorded on a camera sensory array is approximately equal to the convolution of the original image (middle panel) and the motion path of the camera (right panel). Each pixel in the blurry image is the weighted sum of pixels in the original image along the camera motion trajectory. De-blurring an image requires identifying the original image and the motion path from a single blurry image (reproduced from Fergus *et al.* (2006)).

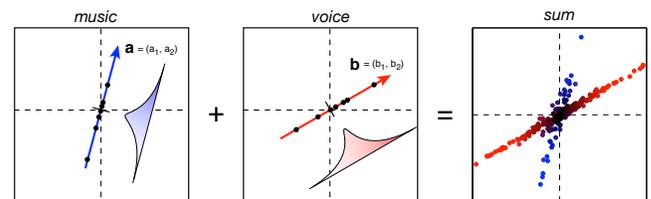


FIG. 3 Example data from the cocktail party problem. Amplitudes recorded simultaneously in microphones 1 and 2 are plotted on the x and y axes, respectively. In the left panel all samples arising from the music lie on the vector $\mathbf{a} = (a_1, a_2)$ reflecting the proximity of the music to microphones 1 and 2, respectively. Likewise, the middle panel depicts all data points arising from solely the voice. The right panel depicts recording of the linear sum both sources. To highlight the contribution of the voice and the music we color each recorded sample by the relative contribution of each each source.

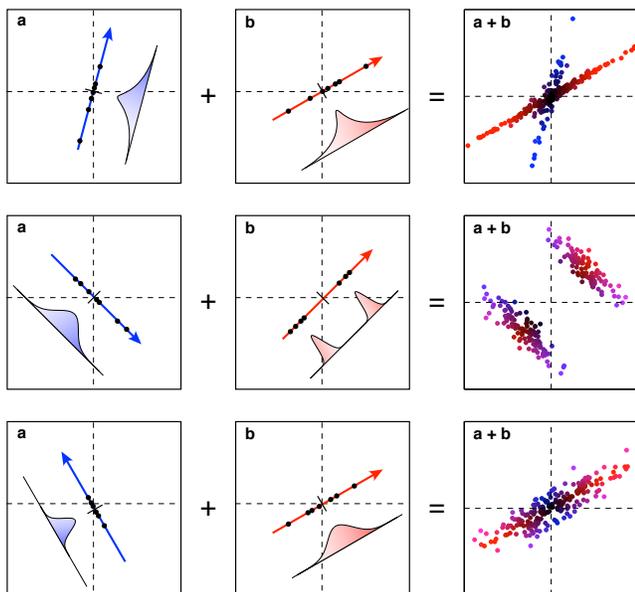


FIG. 4 Examples of linear mixed data. Consider two sources \mathbf{a} and \mathbf{b} represented as blue and red vectors ($n = 2$ dimensions). Each source has a direction represented by the vector termed the *independent components* and a magnitude whose amplitude varies (randomly) according to some distribution. The sum of the independent components weighted by each sample’s magnitude $\mathbf{a} + \mathbf{b}$ produces a data point whose color represents the sum of the magnitudes of the sources which contributed to the data point. *Top row.* The combination of two sharply peaked sources produces an X formation. *Middle row.* Same as top row but the first and second sources are unimodal and bimodal Gaussian distributed (adapted from A. Gretton). *Bottom row.* Same as top but both sources are Gaussian distributed.

III. EXAMPLES ON LINEAR MIXED SIGNALS

Let us imagine data generated from the cocktail party problem in Figure 1. Microphone 1 records the music and voice with proximities a_1 and b_1 , respectively. We plot on the x and y axes the recorded values for microphones 1 and 2, respectively in Figure 3.

What happens if music alone is played? What would the data look like? If music alone is played without a voice, then any data point would lie along the vector (a_1, a_2) . Why? The reason is that (a_1, a_2) measures the proximity of the music to microphones 1 and 2, respectively. Since the music resides at a fixed position with respect to the microphones, each recorded data point must be yoked to (a_1, a_2) as the volume of the music varies over time. Finally, all of the music volumes across time are sampled from some underlying distribution. This behavior is depicted the left panel of Figure 3.

The same experiment could be applied to the data arising from only the voice. This is depicted in the middle panel of Figure 3. Note that the bases associated with each audio source (a_1, a_2) and (b_1, b_2) merely reflect the proximity of each audio source to the microphones. Therefore, the bases need not

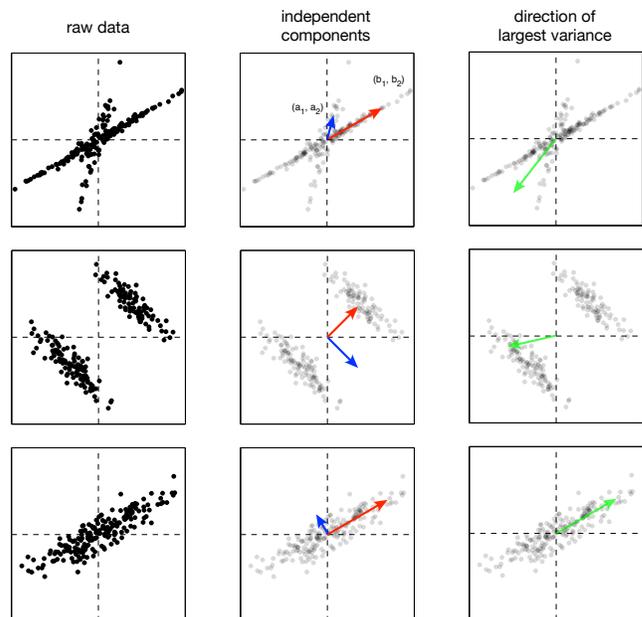


FIG. 5 Analysis of linear mixed signals. *Left column.* Data is reproduced from Figure 4. *Middle column.* Data is replotted superimposed on the basis of the underlying sources, i.e. the independent components (red and blue arrows) in the color corresponding to the basis used to generate the data in Figure 4. *Right column.* Data is replotted superimposed on the direction of largest variance (green arrows). Note that the direction of largest variance might or might not correspond to the independent components.

be orthogonal on the (x, y) plot. Each basis vector labeled $\mathbf{a} = (a_1, a_2)$ and $\mathbf{b} = (b_1, b_2)$, respectively, is termed the independent component (IC) of the data.

Finally, If both the music and voice are played together, the simultaneous measurement in microphones 1 and 2 would be the vector sum of the samples drawn from each basis (again, because sound adds linearly). The composite recording is depicted in the top right panel of Figure 4, where each data point is colored according to the relative contribution of the music and voice (red and blue, respectively).

The middle and bottom rows of Figure 4 depict two different examples of linear mixed data to give a sense of the diversity of BSS problem. In the middle row, the two ICs are orthogonal but the distributions along each axis are dramatically different – namely unimodal and bimodal distributions, respectively. The resulting distribution appears as two oriented lobes. The bottom row depicts two orthogonal ICs, each Gaussian distributed. The resulting distribution is also Gaussian distributed.

In the situation confronted by an experimentalist, the observed data contain no “color labeling” and all analyses must be performed solely on the linear sum (Figure 5, left column). Examining the black points in each example, an experimenter might wish to select a “piece” of the data – for instance, one arm of the X, one mode of the bimodal data or one direction of the Gaussian distribution in each panel respectively. Be-

cause this author judiciously constructed these examples, it is deceptively easy to filter out the distracting source and extract the desired source from each example (e.g. music or voice). Note that these examples are two-dimensional and although we have the benefit of *seeing* the answer, in higher dimensions visualization becomes difficult and a solution might not be as salient.

Because a solution is salient in these examples, it is worthwhile to codify the solution in order to understand what we are aiming for in more complex data sets. The solution to each example is to identify the basis underlying each source, i.e. the independent components. This is depicted by the red and blue arrows in the middle column of Figure 5. The independent components (ICs) might not be orthogonal (top row). Furthermore, the ICs often do not correspond with the direction of maximal variance (top and middle rows). Note that the latter points imply that any technique simply examining the variance would fail to find the ICs (although quizzically the variance does seem to identify one IC in the bottom example). Understanding these issues and how to recover these bases are the heart of this tutorial.

The goal of this tutorial is to build a mathematical solution to the intuition highlighted in all of these examples. We term the solution ICA. During the course of the manuscript, we will understand how and why the types of distributions play an intimate role in this problem and understand how to recover the independent components of most any data set.

IV. SETUP

Here is the framework. We record some multi-dimensional data \mathbf{x} . We posit that each sample is a random draw from an unknown distribution $P(\mathbf{x})$ (e.g. black points in Figure 5). To keep things interesting, we consider \mathbf{x} to contain more than one dimension.³

We assume that there exists some underlying sources \mathbf{s} where each source s_i is statistically independent – the observation of each source s_i is independent of all other sources s_j (where $i \neq j$). For instance, in the cocktail party problem the amplitude of the voice s_1 is independent of the amplitude of the music s_2 at each moment of time.

The key assumption behind ICA is that the observed data \mathbf{x} is a *linear* mixture of the underlying sources

$$\mathbf{x} = \mathbf{A}\mathbf{s} \quad (1)$$

where \mathbf{A} is some unknown invertible, square matrix that mixes the components of the sources. In the example from Figure 1

³ A word about notation. Bold lower case letters \mathbf{x} are column vectors whose i^{th} element is x_i , while bold upper case letters \mathbf{A} are matrices. The probability of observing a random variable Y is formally $P_Y(Y=y)$ but will be abbreviated $P(y)$ for simplicity.

$\mathbf{A} = \begin{bmatrix} a_1 & b_1 \\ a_2 & b_2 \end{bmatrix}$. The goal of ICA is to find the *mixing* matrix \mathbf{A} (more specifically, the inverse of \mathbf{A}) in order to recover the original signals \mathbf{s} from the observed data \mathbf{x} .

We will construct a new matrix \mathbf{W} such that the linear transformed data is an estimate of the underlying sources,

$$\hat{\mathbf{s}} = \mathbf{W}\mathbf{x} \quad (2)$$

In this setting the goal of ICA is to find an *unmixing* matrix \mathbf{W} that is an approximation of \mathbf{A}^{-1} so that $\hat{\mathbf{s}} \approx \mathbf{s}$.

On the face of it, this might appear an impossible problem: find two unknowns \mathbf{A} and \mathbf{s} by only observing their matrix product \mathbf{x} (i.e. black points in Figure 5). Intuitively, this is akin to solving for a and b by only observing $c = a \times b$. Mathematically, one might call this an *under-constrained* problem because the number of unknowns exceed the number of observations.

This intuition is precisely what makes ICA a challenging problem. What I hope to convey to you is that by examining the statistics of the observed data \mathbf{x} , we can find a solution for this problem. This prescription is the essence of ICA.

V. A STRATEGY FOR SOLVING ICA

Divide-and-conquer provides a strategy to solve this problem. Rather than trying to solve for \mathbf{s} and \mathbf{A} simultaneously, we focus on finding \mathbf{A} . Furthermore, rather than trying to solve for \mathbf{A} all at once, we solve for \mathbf{A} in a piece-meal fashion by cutting up \mathbf{A} into simpler and more manageable parts.

The *singular value decomposition* (SVD) is a linear algebra technique that provides a method for dividing \mathbf{A} into several simpler pieces. For any matrix SVD states that

$$\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T.$$

Any matrix is decomposed into three “simpler” linear operations: a rotation \mathbf{V} , a stretch along the axes $\mathbf{\Sigma}$, and a second rotation \mathbf{U} . Each matrix in the SVD is “simpler” because each matrix contains fewer parameters to infer and each matrix is trivial to invert: \mathbf{U} and \mathbf{V} are rotation matrices (or orthogonal matrices) and $\mathbf{\Sigma}$ is a diagonal matrix with real, non-negative values. Figure 6 provides a familiar graphical depiction of SVD.

We estimate \mathbf{A} and its inverse \mathbf{W} by recovering each piece of the decomposition individually:

$$\mathbf{W} = \mathbf{A}^{-1} = \mathbf{V}\mathbf{\Sigma}^{-1}\mathbf{U}^T \quad (3)$$

This equation exploits the fact that the inverse of a rotation matrix is its transpose, (i.e. $\mathbf{V}^{-1} = \mathbf{V}^T$, see Appendix). Furthermore, because \mathbf{A} is invertible by assumption, $\mathbf{\Sigma}^{-1}$ exists and is well defined.

The tutorial will proceed by solving for the unmixing matrix \mathbf{W} in two successive stages:

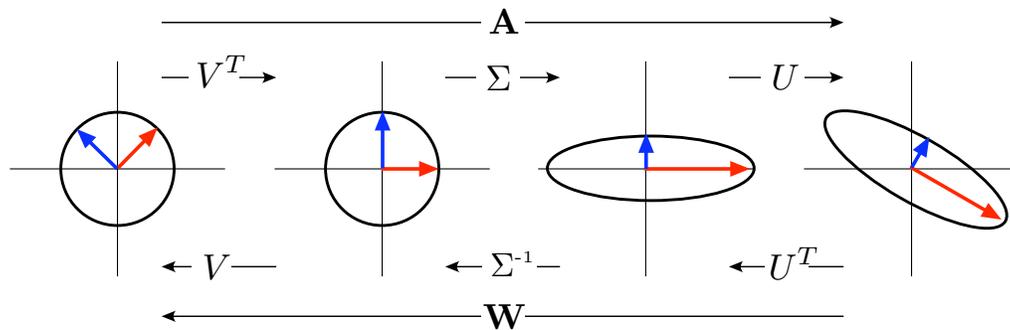


FIG. 6 Graphical depiction of the singular value decomposition (SVD) of matrix $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$ assuming \mathbf{A} is invertible. \mathbf{V} and \mathbf{U} are rotation matrices and $\mathbf{\Sigma}$ is a diagonal matrix. Red and blue arrows are vectors that correspond to the columns of matrix \mathbf{V} (i.e. the basis of the row space of \mathbf{A}). Note how the basis rotates, stretches and rotates during each successive operation. The composition of all three matrix operations is equal to the operation performed by \mathbf{A} . The inverse of matrix \mathbf{A} defined as $\mathbf{W} = \mathbf{V}\mathbf{\Sigma}^{-1}\mathbf{U}^T$ performs each linear operation in the reverse order. Diagram adapted from Strang (1988).

1. Examine the covariance of the data \mathbf{x} in order to calculate \mathbf{U} and $\mathbf{\Sigma}$. (Sections VI and VII)
2. Return to the assumption of independence of \mathbf{s} to solve for \mathbf{V} . (Sections VIII and IX)

Finally, we present a complete solution to ICA, consider the limits and failures of this technique, and offer concrete suggestions for interpreting results from ICA.

At first encounter the divide-and-conquer strategy might appear weakly motivated but in retrospect, dividing the problem into these parts provides a natural strategy that mirrors the structure of correlations in any data set. We will emphasize these larger points throughout the tutorial to provide the reader a framework for thinking about data in general.

VI. EXAMINING THE COVARIANCE OF THE DATA

The goal of this section is to explore the covariance of the data given our assumptions, and in the process recover two of the three matrix operations constituting \mathbf{W} (Equation 3). The covariance of the data provides an appropriate starting point because the covariance matrix measures all correlations that can be captured by a linear model.⁴

As a reminder, the covariance is the expected value of the outer product of individual data points $\langle \mathbf{x}\mathbf{x}^T \rangle$. In order to recover $\mathbf{\Sigma}$ and \mathbf{U} we make one additional assumption in a seemingly unmotivated fashion: assume that the covariance of the sources \mathbf{s} is *whitened*, or equivalently $\langle \mathbf{s}\mathbf{s}^T \rangle = \mathbf{I}$, where \mathbf{I} is

the identity matrix. We discuss what this assumption means in the following section, but for now we make this assumption blindly and will see what this implies about the observed data \mathbf{x} .

The covariance of the data can be expressed in terms of the underlying sources by plugging in the linear mixture model (Equation 1):

$$\begin{aligned} \langle \mathbf{x}\mathbf{x}^T \rangle &= \langle (\mathbf{A}\mathbf{s})(\mathbf{A}\mathbf{s})^T \rangle \\ &= \langle (\mathbf{U}\mathbf{\Sigma}\mathbf{V}^T\mathbf{s})(\mathbf{U}\mathbf{\Sigma}\mathbf{V}^T\mathbf{s})^T \rangle \\ &= \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T \langle \mathbf{s}\mathbf{s}^T \rangle \mathbf{V}\mathbf{\Sigma}\mathbf{U}^T \end{aligned}$$

We exploit our assumption about \mathbf{s} (i.e. $\langle \mathbf{s}\mathbf{s}^T \rangle = \mathbf{I}$) and a property of orthogonal matrices (i.e. $\mathbf{V}^T = \mathbf{V}^{-1}$) to arrive at a final expression

$$\langle \mathbf{x}\mathbf{x}^T \rangle = \mathbf{U}\mathbf{\Sigma}^2\mathbf{U}^T. \quad (4)$$

By our shrewd choice of assumption, note that the covariance of the data is independent of sources \mathbf{s} as well as \mathbf{V} !

What makes Equation 4 extra-special is that it expresses the covariance of the data in terms of a diagonal matrix $\mathbf{\Sigma}^2$ sandwiched between two orthogonal matrices \mathbf{U} . Hopefully, the form of Equation 4 looks familiar to students of linear algebra, but let us make this form explicit.

As an aside, linear algebra tells us that any symmetric matrix (including a covariance matrix) is orthogonally diagonalized by their eigenvectors.⁵ Consider a matrix \mathbf{E} whose columns are the eigenvectors of the covariance of \mathbf{x} . We can prove that

$$\langle \mathbf{x}\mathbf{x}^T \rangle = \mathbf{E}\mathbf{D}\mathbf{E}^T \quad (5)$$

⁴ More specifically, the covariance matrix is a square symmetric matrix where the ij^{th} value is the covariance between x_i and x_j . The ij^{th} term measures all second-order correlations that can be captured by a linear model between x_i and x_j . For simplicity we assume that the mean of the data is zero or that the mean has been subtracted off for each dimension.

⁵ *Orthogonal diagonalization* is a standard decomposition in linear algebra. The term refers to the operation of converting an arbitrary matrix \mathbf{A} into a diagonal matrix by multiplying by an orthogonal basis. Orthogonal diagonalization is achieved by an orthogonal basis of eigenvectors stacked as columns and a diagonal matrix composed of eigenvalues.

where \mathbf{D} is a diagonal matrix of associated eigenvalues (see Appendix). The eigenvectors of the covariance of the data form an orthonormal basis meaning that \mathbf{E} is an orthogonal matrix.

Compare Equations 4 and 5. Both equations state that the covariance of the data can be diagonalized by an orthogonal matrix. Equation 4 provides a decomposition based on the underlying assumption of ICA. Equation 5 provides a decomposition based purely on the properties of symmetric matrices.

Diagonalizing a symmetric matrix using its eigenvectors is a *unique* solution up to a permutation, i.e. no other basis can diagonalize a symmetric matrix. Therefore, if our assumptions behind ICA are correct, then we have identified a partial solution to \mathbf{A} : \mathbf{U} is a matrix of the stacked eigenvectors of the covariance of the data and Σ is a diagonal matrix with the square root of the associated eigenvalue in the diagonal.

Let us summarize our current state. We are constructing a new matrix \mathbf{W} via Equation 3. We have identified the latter two matrices such that

$$\mathbf{W} = \mathbf{V}\mathbf{D}^{-\frac{1}{2}}\mathbf{E}^T.$$

\mathbf{D} and \mathbf{E} are the eigenvalues and eigenvectors of the covariance of the data \mathbf{x} . \mathbf{V} is the sole unknown rotation matrix. Let us now take a moment to interpret our results.

VII. WHITENING, REVISITED

The solution for ICA thus far performs an operation familiar to signal processing termed *whitening*. Whitening is an operation that removes all linear dependencies in a data set (i.e. second-order correlations) and normalizes the variance along all dimensions. Colloquially, this operation is termed *sphereing* the data as intuitively, whitening maps the data into a spherically symmetric distribution.

This intuition is demonstrated by the two operations $\mathbf{D}^{-\frac{1}{2}}\mathbf{E}^T$ depicted in Figure 7. In the first step the data is rotated to align the eigenvectors of the covariance along the cartesian basis. Multiplying by \mathbf{E}^T performs this rotation in order to decorrelate the data, i.e. remove linear dependencies. Mathematically, decorrelation means that the covariance of the transformed data is diagonalized. In our case,

$$\langle (\mathbf{E}^T \mathbf{x}) (\mathbf{E}^T \mathbf{x})^T \rangle = \mathbf{D}$$

where \mathbf{D} is a diagonal matrix of the eigenvalues.⁶ Each diagonal entry in \mathbf{D} is an eigenvalue of the covariance of the data and measures the variance along each dimension.

This operation has a familiar name - *principal component analysis* (PCA). The eigenvectors of the covariance of the data

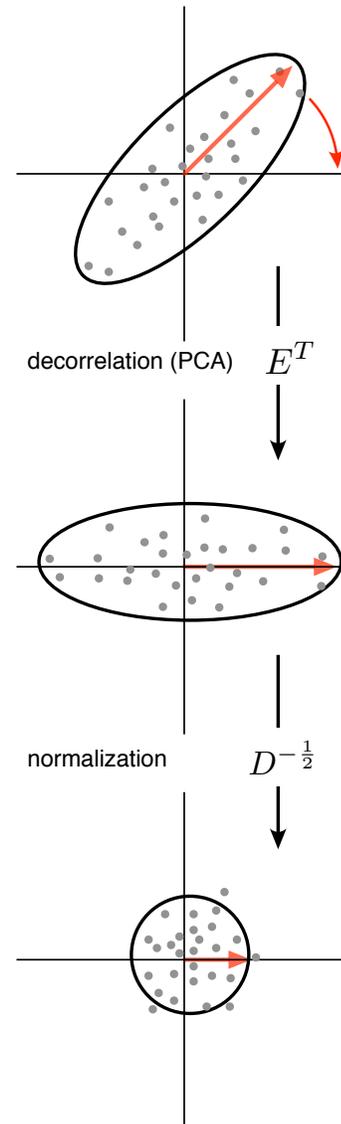


FIG. 7 Whitening a data set can be represented as a series of two linear operations. Data is projected on the principal components, $\mathbf{E}^T \mathbf{x}$. Each axis is then scaled so that every direction has unit variance, $\mathbf{D}^{-\frac{1}{2}} \mathbf{E}^T \mathbf{x}$. The red arrow indicates the transformation of the eigenvector with the largest variance.

from Equation 5,

$$\begin{aligned} \langle \mathbf{x}\mathbf{x}^T \rangle &= \mathbf{E}\mathbf{D}\mathbf{E}^T \\ \mathbf{E}^T \langle \mathbf{x}\mathbf{x}^T \rangle \mathbf{E} &= \mathbf{D} \\ \langle \mathbf{E}^T \mathbf{x}\mathbf{x}^T \mathbf{E} \rangle &= \mathbf{D} \\ \langle (\mathbf{E}^T \mathbf{x}) (\mathbf{E}^T \mathbf{x})^T \rangle &= \mathbf{D} \end{aligned}$$

We have exploited the property that \mathbf{E} is an orthogonal matrix as well as the linearity of expectations.

⁶ This equation directly results from orthogonal diagonalization. Starting

are termed the principal components of the data. Projecting a data set onto the principal components removes linear correlations and provides a strategy for dimensional reduction (by selectively removing dimensions with low variance).

The second operation normalizes the variance in each dimension by multiplying with $\mathbf{D}^{-\frac{1}{2}}$ (Figure 7). Intuitively, normalization ensures that all dimensions are expressed in standard units. No preferred directions exist and the data is rotationally symmetric – much like a sphere.

In our problem whitening simplifies the ICA problem down to finding a single rotation matrix \mathbf{V} . Let us make this simplification explicit by defining

$$\mathbf{x}_w = (\mathbf{D}^{-\frac{1}{2}} \mathbf{E}^T) \mathbf{x}$$

where \mathbf{x}_w is the whitened version of the observed data such that $\langle \mathbf{x}_w \mathbf{x}_w^T \rangle = \mathbf{I}$. Substituting the above equation into Equations 2 and 3 simplifies ICA down to solving $\hat{\mathbf{s}} = \mathbf{V} \mathbf{x}_w$. Note that the problem reduces to finding a rotation matrix \mathbf{V} .

The simplified form of ICA provides additional insight into the structure of the recovered sources $\hat{\mathbf{s}}$. Figure 7 highlights that whitened data \mathbf{x}_w is rotationally symmetric, therefore rotated whitened data $\hat{\mathbf{s}}$ must likewise be whitened (remember $\hat{\mathbf{s}} = \mathbf{V} \mathbf{x}_w$). This is consistent with our assumptions about \mathbf{s} in Section VI. Note that this implies that there exists multiple whitening filters including $\mathbf{D}^{-\frac{1}{2}} \mathbf{E}^T$ and $\mathbf{V}(\mathbf{D}^{-\frac{1}{2}} \mathbf{E}^T)$

VIII. THE STATISTICS OF INDEPENDENCE

The goal of ICA is to find the linear transformation \mathbf{W} that recovers the linear mixed sources \mathbf{s} from the data. By assumption $\mathbf{x} = \mathbf{A} \mathbf{s}$ where \mathbf{x} is the data and both \mathbf{A} and \mathbf{s} are unknown. Exploiting the decomposition of \mathbf{A} and whitening the data has reduced the problem to finding a rotation matrix \mathbf{V} such that $\hat{\mathbf{s}} = \mathbf{V} \mathbf{x}_w$. We now need to exploit the statistics of independence to identify \mathbf{V} .

Remember that the covariance matrix measures the linear dependence between all pairs of variables based on second-order correlations. Whitening the data removed all second-order correlations, hence discerning the last rotation matrix \mathbf{V} requires examining other measures of dependency.

Statistical independence is the strongest measure of dependency between random variables. It requires that neither second-order correlations nor higher-order correlations exist. In particular, if two random variables a and b are independent, then $P(a, b) = P(a) P(b)$ – i.e. the joint probability factorizes. In the context of ICA, we assume that all sources are statistically independent, thus

$$P(\mathbf{s}) = \prod_i P(s_i).$$

This implies that the joint distribution of sources $P(\mathbf{s})$ is a special family of distributions termed a *factorial distribution* be-

cause the joint distribution is the product of the distribution of each source $P(s_i)$.

The problem of ICA searches for the rotation \mathbf{V} such that $\hat{\mathbf{s}}$ is statistically independent $P(\hat{\mathbf{s}}) = \prod_i P(\hat{s}_i)$. Because all second-order correlations have been removed, the unknown matrix \mathbf{V} must instead remove all higher-order correlations. Removing all higher-order correlations with a single rotation \mathbf{V} is a tall order, but if the model $\mathbf{x} = \mathbf{A} \mathbf{s}$ is correct, then this is achievable and $\hat{\mathbf{s}}$ will be statistically independent.

We therefore require a function (termed a *contrast function*) that measures the amount of higher-order correlations – or equivalently, measures how close the estimated sources $\hat{\mathbf{s}}$ are to statistical independence.

IX. A SOLUTION USING INFORMATION THEORY

The entire goal of this section is to find the last rotation \mathbf{V} such that the estimate of $\hat{\mathbf{s}}$ is statistically independent. To solve this problem, we resort to a branch of mathematics called *information theory*. The mathematics of information theory is quite distinct. I would suggest that those new to this material skim this section in their first reading.

Many contrast functions exist for measuring statistical independence of $\hat{\mathbf{s}}$. Examples of contrast functions include rigorous statistical measures of correlations, approximations to said measures (e.g. see code in Appendix B) and clever, ad-hoc guesses.

Because this tutorial is focused on presenting the foundational ideas behind ICA, we focus on a natural measure from information theory to judge how close a distribution is to statistical independence. The *mutual information* measures the departure of two variables from statistical independence. The *multi-information*, a generalization of mutual information, measures the statistical dependence between multiple variables:

$$I(\mathbf{y}) = \int P(\mathbf{y}) \log_2 \frac{P(\mathbf{y})}{\prod_i P(y_i)} d\mathbf{y}$$

It is a non-negative quantity that reaches a minimum of zero if and only if all variables are statistically independent. For example, if $P(\mathbf{y}) = \prod_i P(y_i)$, then $\log(1) = 0$ and $I(\mathbf{y}) = 0$.

The goal of ICA can now be stated succinctly. Find a rotation matrix \mathbf{V} such that $I(\hat{\mathbf{s}}) = 0$ where $\hat{\mathbf{s}} = \mathbf{V} \mathbf{x}_w$. If we find such a rotation, then $\hat{\mathbf{s}}$ is statistically independent. Furthermore, $\mathbf{W} = \mathbf{V} \mathbf{D}^{-\frac{1}{2}} \mathbf{E}^T$ is the solution to ICA and we can use this matrix to estimate the underlying sources.

The multi-information is minimized when $\hat{\mathbf{s}}$ is statistically independent, therefore the goal of ICA is to minimize the multi-information until it reaches zero. Reconsider the data from the middle example in Figure 5. \mathbf{V} is a rotation matrix and in two dimensions \mathbf{V} has the form

$$\mathbf{V} = \begin{bmatrix} \cos(\theta) & -\sin(\theta) \\ \sin(\theta) & \cos(\theta) \end{bmatrix}.$$

The rotation angle θ is the only free variable. We can calculate the multi-information of $\hat{\mathbf{s}}$ for all θ in Figure 8. Examining the plot we notice that indeed the multi-information is zero when the rotation matrix is 45° . This means that the recovered distributions are statistically independent. The minimum of the multi-information is zero implying that \mathbf{V} is a two-dimensional 45° rotation matrix.

On the surface the optimization of multi-information might appear abstract. This procedure though can be visualized (and validated) by plotting the recovered sources $\hat{\mathbf{s}}$ using different rotations (Figure 8, bottom). At a 45° rotation the recovered sources $\hat{\mathbf{s}}$ along the x-axis and y-axis are the original bimodal and unimodal gaussian distributions, respectively (Figure 8, bottom left). Note though that the optimal rotation is not unique as adding integer multiples of 90° also minimizes the multi-information!

Minimizing the multi-information is difficult in practice but can be simplified. A simplified form of the optimization bares important relationships with other interpretations of ICA.

The multi-information is a function of the entropy $H[\cdot]$ of a distribution. The entropy $H[\mathbf{y}] = -\int P(\mathbf{y}) \log_2 P(\mathbf{y}) d\mathbf{y}$ measures the amount of uncertainty about a distribution $P(\mathbf{y})$. The multi-information is the difference between the sum of entropies of the marginal distributions and the entropy of the joint distribution, i.e. $I(\mathbf{y}) = \sum_i H[y_i] - H[\mathbf{y}]$. Therefore, the multi-information of $\hat{\mathbf{s}}$ is

$$\begin{aligned} I(\hat{\mathbf{s}}) &= \sum_i H[(\mathbf{V}\mathbf{x}_w)_i] - H[\mathbf{V}\mathbf{x}_w] \\ &= \sum_i H[(\mathbf{V}\mathbf{x}_w)_i] - (H[\mathbf{x}_w] + \log_2 |\mathbf{V}|) \end{aligned} \quad (6)$$

where $(\mathbf{V}\mathbf{x}_w)_i$ is the i^{th} element of $\hat{\mathbf{s}}$ and we have employed an expression relating the entropy of a probability density under a linear transformation.⁷ The determinant of a rotation matrix is 1 so the last term is zero, i.e. $\log_2 |\mathbf{V}| = 0$.

The optimization is simplified further by recognizing that we are only interested in finding the rotation matrix (and not the value of the multi-information). The term $H[\mathbf{x}_w]$ is a constant and independent of \mathbf{V} , so it can be dropped:

$$\mathbf{V} = \arg \min_{\mathbf{V}} \sum_i H[(\mathbf{V}\mathbf{x}_w)_i] \quad (6)$$

The optimization has simplified to finding a rotation matrix that minimizes the sum of the marginal entropies of $\hat{\mathbf{s}}$. The rotation matrix \mathbf{V} that solves Equation 6 maximizes the statistical independence of $\hat{\mathbf{s}}$.

A few important connections should be mentioned at this point. First, calculating the entropy from a finite data set is difficult and should be approached with abundant caution.⁸

⁷ For a linear transformation \mathbf{B} and a random variable \mathbf{x} , $H[\mathbf{B}\mathbf{x}] = H[\mathbf{x}] + \log_2 |\mathbf{B}|$.

⁸ The entropy is a function of an unknown distribution $P(\mathbf{y})$. Instead any

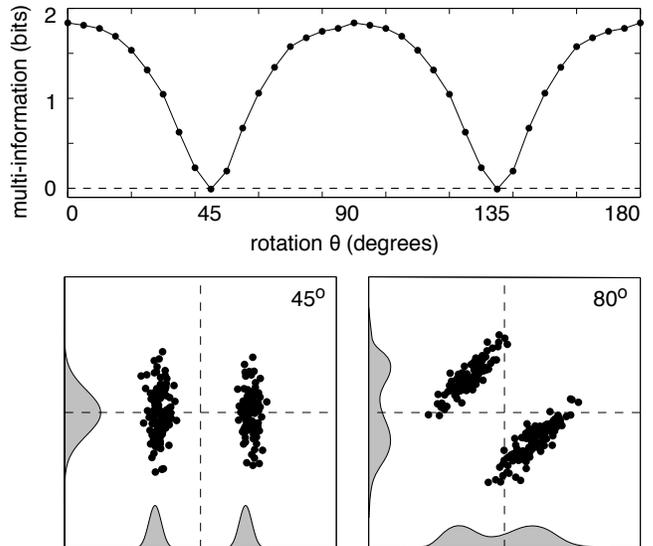


FIG. 8 Optimization to recover statistical independence for the middle example from Figure 5 (see Appendix B). Rotation matrix \mathbf{V} contains one free variable, the rotation angle θ . The multi-information of $\hat{\mathbf{s}}$ is plotted versus the rotation angle θ (top). Recovered sources $\hat{\mathbf{s}}$ are plotted for two estimates of $\mathbf{W} = \mathbf{V}\mathbf{D}^{-\frac{1}{2}}\mathbf{E}^T$ where \mathbf{V} is a rotation matrix with a 45° and 80° rotation, respectively (bottom). Grey curves are marginal distributions along each dimension (bottom, inset). Note that a rotation of 45° degrees recovers a bimodal and unimodal gaussian distribution along each axis.

Thus, many ICA optimization strategies focus on approximations to Equation 6 (but see Lee *et al.* (2003)).

Second, the form of Equation 6 has multiple interpretations reflecting distinct but equivalent interpretations of ICA. In particular, a solution to Equation 6 finds the rotation that maximizes the “non-Gaussianity” of the transformed data.⁹ Likewise, Equation 6 is equivalent to finding the rotation that maximizes the log-likelihood of the observed data under the assumption that the data arose from a statistically independent distribution. Both interpretations provide starting points that other authors have employed to derive Equation 6.

In summary, we have identified an optimization (Equation 6) that permits us to estimate \mathbf{V} and in turn, reconstruct the original statistically independent source signals $\hat{\mathbf{s}} = \mathbf{W}\mathbf{x}$. The columns of \mathbf{W}^{-1} are the independent components of the data.

practitioner solely has access to a finite number of samples from this distribution. This distinction is of paramount importance because any estimate of entropy based on finite samples of data can be severely biased. This technical difficulty has motivated new methods for estimating entropy which minimize such biases (see code in Appendix).

⁹ The deviation of a probability distribution from a Gaussian is commonly measured by the *negentropy*. The negentropy is defined as the *Kullback-Leibler divergence* of a distribution from a Gaussian distribution with equal variance. Negentropy is equal to a constant minus Equation 6 and maximizing the sum of the marginal negentropy (or non-Gaussianity) is equivalent to minimizing the sum of the marginal entropy.

Quick Summary of ICA

1. Subtract off the mean of the data in each dimension.
2. Whiten the data by calculating the eigenvectors of the covariance of the data.
3. Identify final rotation matrix that optimizes statistical independence (Equation 6).

FIG. 9 Summary of steps behind ICA. The first two steps have analytical solutions but the final step must be optimized numerically. See Appendix B for example code.

X. DISCUSSION

The goal of ICA is to recover linear mixed signals without knowledge of the linear mixture. ICA recovers the linear mixture and the underlying sources by assuming that the signals are statistically independent. Although a recent addition to the modern toolbox, ICA has found increasing popularity in signal processing and machine learning for filtering signals and performing dimensional reduction (for summary, see Figure 9; for example code, see Appendix B).

Identifying the linear mixture, i.e. the basis of independent components, provides a means for selecting (or deleting) individual sources of interest. For instance, in the cocktail party example in Figure 1, each recovered independent component is a column of \mathbf{W}^{-1} and corresponds to the filter selecting the voice or the background music from the pair of microphones. Multiplying the data \mathbf{x} by an individual row of the unmixing matrix \mathbf{W} recovers an estimate of the voice or the music respectively. Note that if one records some new data from the microphones, then the unmixing matrix could be applied to the new data as well.

The success of ICA must be tempered with several computational issues and inherent ambiguities in the solution. In the computation of the unmixing matrix $\mathbf{W} = \mathbf{V}\Sigma^{-1}\mathbf{U}^T$, the matrices Σ^{-1} and \mathbf{U} are analytically calculated from the data. The matrix \mathbf{V} has no analytical form and must be approximated numerically through an optimization procedure. The optimization is inherently difficult because local minima exist in the objective function (Equation 6) complicating any procedure based on ascending a gradient. In addition, estimating the quantity that must be optimized (the entropy of a distribution) from a finite number of data points is extremely difficult. These two challenges are often approached by approximating the entropy with related statistics (e.g. measures of correlation) that are easier to optimize with finite data sets.

Even when the above issues are addressed, any ICA solution is subject to several inherent ambiguities. These ambiguities exist because the objective is agnostic to these degrees of freedom (Figure 10). First, the labels of each independent

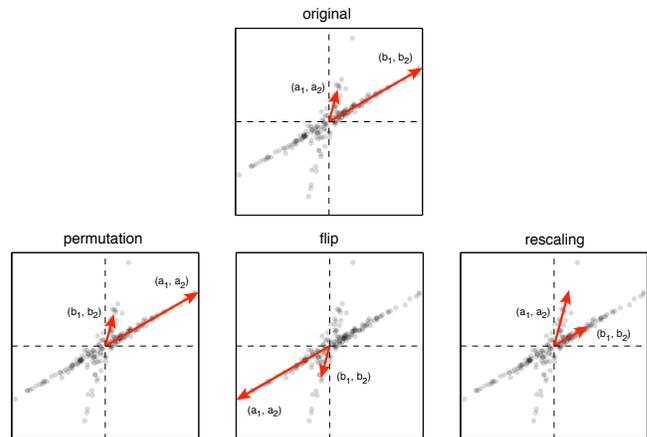


FIG. 10 Ambiguities in the ICA solution. Each example demonstrates a degree of freedom that provides an additional solution for the top panel (reproduction of Figure 5). *Left*. A permutation of the labels of the recovered independent components. *Middle*. A flip of the recovered independent components. *Right*. A rescaling of the recovered independent components.

component can be arbitrarily permuted.¹⁰ Second, any independent component can be flipped across the origin. Third, any independent component can be rescaled with arbitrary length (because one can absorb any rescaling into the inferred sources $\hat{\mathbf{s}}$).¹¹ For most applications these ambiguities are inconsequential, however being cognizant of these issues could prevent an inappropriate interpretation of the results.

Given the challenges and ambiguities in the ICA solution highlighted above, one might question what benefit exists in analyzing a data set with ICA. The primary advantage of ICA is that it provides a fairly unrestricted linear basis for representing a data set. This advantage is often highlighted by contrasting the results of ICA with PCA. PCA identifies an orthogonal linear basis which maximizes the variance of the data. ICA is not restricted to an orthogonal basis because statistical independence makes no such requirement. This is most apparent in the top example of Figure 5 where the directions of interest are quite salient, but maximizing the explained variance fails to identify these directions. In terms of real world data, this lack of restriction means that independent sources need not be orthogonal but merely linear independent.

¹⁰ Since \mathbf{W} undoes the linear transformation of \mathbf{A} , we might infer that $\mathbf{W}\mathbf{A} = \mathbf{I}$, where \mathbf{I} is an identity matrix. This is overly restrictive, however any permutation of the identity matrix is also a valid ICA solution, e.g.

$$\mathbf{W}\mathbf{A} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \text{ or } \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}.$$

¹¹ We resolve the last ambiguity by selecting an arbitrary scale for the recovered sources, in particular a scale that simplifies the computation. For instance, earlier in the text we assumed that the covariance of the source is white $\langle \mathbf{s}\mathbf{s}^T \rangle = \mathbf{I}$. This assumption is a trick that accomplishes two simultaneous goals. First, a variance of 1 in each dimension provides a fixed scale for \mathbf{s} . Second, a white covariance simplifies the evaluation of the data covariance by removing any dependence on \mathbf{s} and \mathbf{V} (Equation 4).

These statements are only part of the story as evident in the middle and bottom examples of Figure 5. In both examples the underlying sources are orthogonal. If orthogonality were the only relevant factor, PCA should have recovered the underlying sources in the middle example. Why does ICA uniquely recover the sources in the middle example? The answer lies in the fact that PCA can identify the underlying sources only when data is distributed appropriately, such as a Gaussian distribution. This statement will be justified in the following section along with an analogous statement for ICA. Importantly, this discussion will lead to conditions under which ICA will succeed or fail to recover the underlying sources.

A. What data works and why

Predicting the success of any dimensional reduction technique ultimately requires a discussion about how the data was generated. The discussion draws heavily on the topic of whitening in Section VII, although the mathematics are slightly more involved.

Let us begin by examining our assumptions closer. We assumed that the underlying sources are statistically independent, i.e. $P(\mathbf{s}) = \prod_i P(s_i)$. As mentioned in Section VII, distributions of this form are termed factorial distributions. Furthermore, the linear mixture model $\mathbf{x} = \mathbf{A}\mathbf{s}$ provides a concrete prescription for how the observed data \mathbf{x} was generated. By assumption ICA expects the data to arise from an (invertible) linear transformation applied to a factorial distribution. Hence, ICA expects and will work optimally on data arising from *linearly transformed factorial distributions* – applying ICA to any other type of data is a gamble with no guarantees.¹²

The case of Gaussian distributed data provides an important didactic point about the application of ICA (Figure 5, bottom example). Why is principal component analysis (PCA) identical to independent component analysis for this data?

As a side bar, remember that the first part of the ICA solution employs PCA to remove second-order correlations (Section VII). The second part of the ICA solution employs the rotation \mathbf{V} to remove higher-order correlations (Section VIII).

Since PCA is equivalent to ICA for Gaussian data, we infer that Gaussian data contains no higher order correlations – but we need not rely on inference to understand this part. If the observed data \mathbf{x} is Gaussian distributed with covariance Σ , ap-

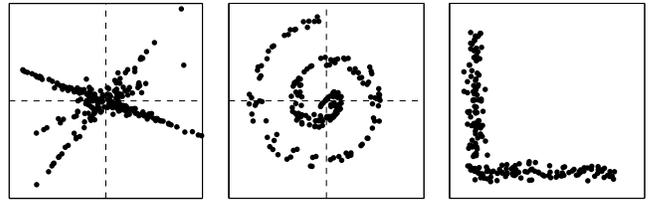


FIG. 11 Challenges to identifying independent components. *Left*: Overcomplete representation. The number of independent sources exceed the number of measurements. *Middle*: Nonlinear manifold. The underlying source lies on a nonlinear manifold. *Right*: Occlusion. The superposition of the two sources is not linear but mutually exclusive.

plying a whitening filter results in whitened data \mathbf{x}_w Gaussian distributed with no covariance structure (i.e. identity covariance matrix). Note that $P(\mathbf{x}_w)$ can be written as a factorial distribution $P(\mathbf{x}_w) = \prod_i \exp(x_{w,i})$. Hence, a whitening filter based on PCA already achieves a factorial distribution and no need exists for computing \mathbf{V} .

This discussion has been theoretical, but it has important implications for the application of ICA on real data. Often, a practitioner of ICA does not know how their data arose and merely tries this technique to see if the results are reasonable. One lesson from this discussion is to ask how prominent higher-order correlations are in the data. If higher-order correlations are small or insignificant, then PCA is sufficient to recover the independent sources of the data.

A second lesson is to ask how independent the recovered sources are. Empirically measuring $I(\hat{\mathbf{s}})$ (or any measure of correlation) provides a useful check to ensure that the sources are indeed independent. In practice, $I(\hat{\mathbf{s}})$ rarely achieves 0 either because of sampling issues (see Appendix) or because of local minima in the optimization. Therefore, the final results must be rigorously checked by calculating multiple measures of correlation or employing bootstrap procedures.

B. Extensions

The tutorial up to this point has offered a foundation for approaching the topic of ICA. A judicious reader might have already recognized several deficiencies in the ICA framework which could readily be improved or extended. The purpose of this section is to highlight a few future directions that a reader would now be prepared to explore.

An obvious shortcoming of ICA can be seen as far back as Figure 1. In retrospect, it seems artificial to stipulate two microphones for separating the two auditory sources. Clearly, a single microphone should suffice for discerning between a potentially large number of auditory sources. Mathematically, this intuition suggests that the number of sources need not match the number of recordings. In particular, the dimension of \mathbf{s} might exceed the dimension of \mathbf{x} . This *overcomplete* problem is a situation in which \mathbf{A} is not a square matrix (Fig-

¹² In general, it is difficult to write down a succinct, closed-form equation for a linear transformed factorial distribution. The Gaussian distribution is an exception to this statement (see below). In general, the marginal distribution of a linear transformed factorial distribution is the linear sum of two independent distributions. The linear sum of two independent distributions is mathematically equal to the *convolution* of the independent distributions with weights given by the coefficients of the independent components.

ure 11, left). Therefore, \mathbf{A} is not an invertible matrix, violating an important assumption behind the discussion of ICA.

If \mathbf{A} is not invertible, then outside information must be employed to solve the problem. Additional information often takes the form of adding a regularization term, or positing a particular distribution for the data. The latter technique is a Bayesian perspective in which one posits that the data arose from a *prior distribution* and the observations are corrupted by some noise process stipulated by a likelihood function. Several versions of ICA have been developed that follow this Bayesian perspective and provide solutions (although limited) to various over-complete situations.

Multiple extensions have further been suggested for ICA based on empirical observations of blind source separation problems in the real world. For instance, some recent extensions have focused on handling data arising from binary or sparse sources, nonlinear interactions (e.g. occlusion) or nonlinear representations. For those who wish to delve into this topic further, a nice starting place are the popular *Infomax* papers (Bell and Sejnowski, 1995, 1997), a special journal issue on ICA (Lee *et al.*, 2003) and an entire text book on this topic (Hyvärinen *et al.*, 2004). Lyu and Simoncelli (2009) provide a tremendous discussion on the statistics of ICA and its relationship to natural signal statistics. Finally, Faivishevsky and Goldberger (2008) provide an elegant implementation mirroring this tutorial's explanations that achieves state-of-the-art performance.

Writing this paper has been an extremely instructional experience for me. I hope that this paper helps to demystify the motivation and results of ICA, and the underlying assumptions behind this important analysis technique. Please send me a note if this has been useful to you as it inspires me to keep writing!

References

- Bell, A., and T. Sejnowski, 1995, *Neural computation* **7**(6), 1129.
 Bell, A., and T. Sejnowski, 1997, *Vision research* **37**(23), 3327.
 Bregman, A., 1994, *Auditory scene analysis: The perceptual organization of sound* (The MIT Press).
 Cardoso, J., 1989, in *Acoustics, Speech, and Signal Processing, 1989. ICASSP-89., 1989 International Conference on* (IEEE), pp. 2109–2112.
 Faivishevsky, L., and J. Goldberger, 2008, *Advances in Neural Information Processing Systems* **11**.
 Fergus, R., B. Singh, A. Hertzmann, S. T. Roweis, and W. Freeman, 2006, *ACM Transactions on Graphics* **25**, 787.
 Hyvärinen, A., J. Karhunen, and E. Oja, 2004, *Independent component analysis*, volume 46 (John Wiley & Sons).
 Lee, T., J. Cardoso, and S. Amari, 2003, *Journal of Machine Learning Research* **4**, 1175.
 Lyu, S., and E. Simoncelli, 2009, *Neural Computation* **21**(6), 1485.
 Strang, G., 1988, *Linear Algebra and Its Applications* (Brooks Cole).
 Victor, J., 2002, *Physical Review E* **66**(5), 051903.

Appendix A: Mathematical Proofs

1. The inverse of an orthogonal matrix is its transpose.

Let \mathbf{A} be an $m \times n$ orthogonal matrix where \mathbf{a}_i is the i^{th} column vector. The ij^{th} element of $\mathbf{A}^T \mathbf{A}$ is

$$(\mathbf{A}^T \mathbf{A})_{ij} = \mathbf{a}_i^T \mathbf{a}_j = \begin{cases} 1 & \text{if } i = j \\ 0 & \text{otherwise} \end{cases}$$

Therefore, because $\mathbf{A}^T \mathbf{A} = \mathbf{I}$, it follows that $\mathbf{A}^{-1} = \mathbf{A}^T$.

2. A symmetric matrix is diagonalized by a matrix of its orthonormal eigenvectors.

Let \mathbf{A} be a square $n \times n$ symmetric matrix with associated eigenvectors $\{\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_n\}$. Let $\mathbf{E} = [\mathbf{e}_1 \ \mathbf{e}_2 \ \dots \ \mathbf{e}_n]$ where the i^{th} column of \mathbf{E} is the eigenvector \mathbf{e}_i . This theorem asserts that there exists a diagonal matrix \mathbf{D} such that $\mathbf{A} = \mathbf{E} \mathbf{D} \mathbf{E}^T$.

This proof is in two parts. In the first part, we see that the any matrix can be orthogonally diagonalized if and only if that matrix's eigenvectors are all linear independent. In the second part of the proof, we see that a symmetric matrix has the special property that all of its eigenvectors are not just linear independent but also orthogonal, thus completing our proof.

In the first part of the proof, let \mathbf{A} be just some matrix, not necessarily symmetric, and let it have independent eigenvectors (i.e. no degeneracy). Furthermore, let $\mathbf{E} = [\mathbf{e}_1 \ \mathbf{e}_2 \ \dots \ \mathbf{e}_n]$ be the matrix of eigenvectors placed in the columns. Let \mathbf{D} be a diagonal matrix where the i^{th} eigenvalue is placed in the i^{th} position.

We will now show that $\mathbf{A} \mathbf{E} = \mathbf{E} \mathbf{D}$. We can examine the columns of the right-hand and left-hand sides of the equation.

$$\begin{aligned} \text{Left hand side : } \mathbf{A} \mathbf{E} &= [\mathbf{A} \mathbf{e}_1 \ \mathbf{A} \mathbf{e}_2 \ \dots \ \mathbf{A} \mathbf{e}_n] \\ \text{Right hand side : } \mathbf{E} \mathbf{D} &= [\lambda_1 \mathbf{e}_1 \ \lambda_2 \mathbf{e}_2 \ \dots \ \lambda_n \mathbf{e}_n] \end{aligned}$$

Evidently, if $\mathbf{A} \mathbf{E} = \mathbf{E} \mathbf{D}$ then $\mathbf{A} \mathbf{e}_i = \lambda_i \mathbf{e}_i$ for all i . This equation is the definition of the eigenvalue equation. Therefore, it must be that $\mathbf{A} \mathbf{E} = \mathbf{E} \mathbf{D}$. A little rearrangement provides $\mathbf{A} = \mathbf{E} \mathbf{D} \mathbf{E}^{-1}$, completing the first part the proof.

For the second part of the proof, we show that a symmetric matrix always has orthogonal eigenvectors. For some symmetric matrix, let λ_1 and λ_2 be distinct eigenvalues for eigenvectors \mathbf{e}_1 and \mathbf{e}_2 .

$$\begin{aligned} \lambda_1 \mathbf{e}_1 \cdot \mathbf{e}_2 &= (\lambda_1 \mathbf{e}_1)^T \mathbf{e}_2 \\ &= (\mathbf{A} \mathbf{e}_1)^T \mathbf{e}_2 \\ &= \mathbf{e}_1^T \mathbf{A}^T \mathbf{e}_2 \\ &= \mathbf{e}_1^T \mathbf{A} \mathbf{e}_2 \\ &= \mathbf{e}_1^T (\lambda_2 \mathbf{e}_2) \\ \lambda_1 \mathbf{e}_1 \cdot \mathbf{e}_2 &= \lambda_2 \mathbf{e}_1 \cdot \mathbf{e}_2 \end{aligned}$$

By the last relation we can equate that $(\lambda_1 - \lambda_2)\mathbf{e}_1 \cdot \mathbf{e}_2 = 0$. Since we have conjectured that the eigenvalues are in fact unique, it must be the case that $\mathbf{e}_1 \cdot \mathbf{e}_2 = 0$. Therefore, the eigenvectors of a symmetric matrix are orthogonal.

Let us back up now to our original postulate that \mathbf{A} is a symmetric matrix. By the second part of the proof, we know that the eigenvectors of \mathbf{A} are all orthonormal (we choose the eigenvectors to be normalized). This means that \mathbf{E} is an orthogonal matrix so by theorem 1, $\mathbf{E}^T = \mathbf{E}^{-1}$ and we can rewrite the final result.

$$\mathbf{A} = \mathbf{E}\mathbf{D}\mathbf{E}^T$$

Thus, a symmetric matrix is diagonalized by a matrix of its eigenvectors.

3. For any zero-mean, random vector \mathbf{x} , the matrix \mathbf{W} whitens \mathbf{x} if and only if $\mathbf{W}^T\mathbf{W} = \langle \mathbf{x}\mathbf{x}^T \rangle^{-1}$.

The definition of whitening \mathbf{x} is that the covariance of $\mathbf{W}\mathbf{x}$ is the identity matrix for some matrix \mathbf{W} . We begin with this definition and observe what this requires of \mathbf{W} .

$$\begin{aligned} \langle (\mathbf{W}\mathbf{x})(\mathbf{W}\mathbf{x})^T \rangle &= \mathbf{I} \\ \mathbf{W} \langle \mathbf{x}\mathbf{x}^T \rangle \mathbf{W}^T &= \mathbf{I} \\ \langle \mathbf{x}\mathbf{x}^T \rangle &= \mathbf{W}^{-1}\mathbf{W}^{-T} \\ \langle \mathbf{x}\mathbf{x}^T \rangle^{-1} &= \mathbf{W}^T\mathbf{W} \end{aligned}$$

The left-hand side is the inverse of the covariance of the data \mathbf{x} , also known as the precision matrix. Note that the matrix \mathbf{W} contains n^2 unknown terms but the equation provides fewer constraints. Multiple solutions for \mathbf{W} must exist in order to whiten the data \mathbf{x} .

Appendix B: Code

This code is written for Matlab 7.4 (R2007a). Several popular packages exist for computing ICA including FastICA and InfoMax which employs several notable algorithmic improvements for speed and performance.

For the purposes of this tutorial, I am including an implementation of a simple ICA algorithm that attempts to find the direction of maximum kurtosis (Cardoso, 1989). The FOBI algorithm is not very robust algorithm although it is quite elegant, simple to code and useful as a didactic tool.

```
function [W, S] = ica(X)
% ICA Perform independent component analysis.
%
% [W, S] = ica(X);
%
% where X = AS and WA = eye(d)
% and [d,n] = size(X) (d dims, n samples).
%
```

```
% Implements FOBI algorithm.
%
% JF Cardoso (1989) "Source separation using
% higher order moments", Proc Intl Conf Acoust
% Speech Signal Process

[d, n] = size(X);

% Subtract off the mean of each dimension.
X = X - repmat(mean(X,2),1,n);

% Calculate the whitening filter.
[E, D] = eig(cov(X'));

% Whiten the data
X_w = sqrtm(pinv(D))*E'*X;

% Calculate the rotation that aligns with the
% directions which maximize fourth-order
% correlations. See reference above.
[V,s,u] = svd((repmat(sum(X_w.*X_w,1),d,1).*X_w)*X_w');

% Compute the inverse of A.
W = V * sqrtm(pinv(D)) * E';

% Recover the original sources.
S = W * X;
```

The next section of code is a script to demonstrate the example from Figure 8.

```
% Parameters of the toy example.
n = 200; % number of samples (quick)
%n = 1000; % number of samples (precise)

sd1 = 0.14; % standard deviation of source 1
sd2 = 0.5; % standard deviation of source 2
angle = 45; % rotation angle

% Generate the bimodal data set.
randn('state',1);
s1 = (sd1 * randn(1,n) + sign(randn(1,n)));
s2 = sd2 * randn(1,n);
S = [s1; s2];

% Generate the mixing matrix to be a rotation.
theta = angle / 360 * 2*pi;
A = [cos(theta) sin(theta); ...
     -sin(theta) cos(theta)];

% Linearly mix the data.
X = A * S;

% Subtract off the mean of each dimension.
X = X - repmat(mean(X,2),1,n);

% Calculate the whitening filter.
[E, D] = eig(cov(X'));

% Whiten the data
X_w = sqrtm(pinv(D))*E'*X;

% Create an array of angles between 0 and 180 deg.
angles = 0:5:180;
thetas = angles / 360 * 2*pi;

% Calculate the multi-information for all angles.
for i=1:length(thetas)
```

```

% Generate a rotation matrix
V = [cos(thetas(i)) sin(thetas(i)); ...
     -sin(thetas(i)) cos(thetas(i))];

% Try to recover the sources using the rotation.
S_est = V * X_w;

% Calculate the multi-information.
I(i) = entropy(S_est(1,:)) + ...
      entropy(S_est(2,:)) - ...
      entropy(S_est);
end

% Plot the multi-information.
figure;
hold on; box on
plot(angles, I, '.k-', 'MarkerSize', 16);
xlabel('angle (degrees)');
ylabel('multi information');
xlim([angles(1) angles(end)]);
plot(xlim, zeros(2,1), '--k');

% Plot the original data with the IC's.
figure;
subplot(1,2,1);
hold on; box on;
plot(X(1,:), X(2,:), '.k', 'MarkerSize', 16);
axis(2*[-1 1 -1 1]); axis square;

% Plot the ICA solution.
[W, S] = ica(X);
subplot(1,2,2);
hold on; box on
plot(S(1,:), S(2,:), '.k', 'MarkerSize', 16);
axis(2*[-2 2 -1 1]); axis square;

```

This function demonstrates a clever binless estimator of entropy which exhibits good statistical properties (Victor, 2002).

```

function H = entropy(X)
% ENTROPY Estimate the entropy using
% a cool binless estimator.

```

```

%
% Note [d,n] = size(X) (d dims, n samples).
%
% J Victor (2002) "Binless strategies for
% estimation of information from neural
% data", Physical Review E.

% Use the machine precision for neighbor
% calculations.
precision = eps;

[d, n] = size(X);

% Calculate the nearest neighbor for each point.
for i=1:n
    % Remove the i'th vector.
    X_temp = X(:, find([1:n] - i));

    % Subtract off the i'th vector from all others.
    X_diff = X_temp - repmat(X(:,i), 1, n-1);

    % Calculate the minimum Euclidean distance.
    lambda(i) = min(sqrt(sum((X_diff).^2, 1)));

    % Ensure the distance is not zero.
    if (lambda(i) < precision)
        lambda(i) = precision;
    end
end

% The "Euler-Mascheroni" constant.
em = 0.5772156649015;

% Calculate area the of an d-dimensional sphere.
area = d * pi^(d/2) / gamma(1 + d/2);

% Calculate an estimate of entropy based on the
% mean nearest neighbor distance using an equation
% from the above citation.
K = log2(area) + log2((n-1) / d) + em / log(2);
H = d * mean(log2(lambda)) + K;

```