

# Data Decomposition: From Independent Component Analysis to Sparse Representations

E. Roussos

MLRG, University of Oxford, Parks Road, Oxford, OX1 3PJ, U.K.

[eroussos@learning-machines.org](mailto:eroussos@learning-machines.org)

## Abstract

This paper provides a unifying review of some recent approaches to decomposing data into sets of components. We start from the classical algebraic method of singular value decomposition and then introduce principal and independent component analysis. The text continues with the main subject of this paper, sparse representation and decomposition, emphasizing its biological plausibility. In this paper emphasis will be given to the geometric perspective, with the mathematics kept to an essential minimum.

## 1 Data modelling: A probabilistic approach

In an exploratory approach to data analysis, it is often useful to consider the observations as generated from a set of latent generators or ‘sources’ via a generally unknown mapping. Our goal is to recover the generators from the observations, an *inverse* problem. This can be often stated as a *data decomposition* problem: the data matrix is decomposed into factors, each one of them representing some salient characteristics of the data. In fact, many well known algorithms, such as singular value decomposition (SVD) and principal component analysis (PCA), independent component analysis (ICA), as well as  $k$ -means and many others can be stated under this formulation, providing a unifying framework for unsupervised learning. Another view is that of the *representation* of data sets in a new coordinate system such that certain properties hold. For example, in PCA we seek a new coordinate system in which the data become linearly uncorrelated. For the noisy overcomplete case, where we have more sources than observations, the problem of reconstructing the sources becomes extremely ill-posed. Solutions to such inverse problems can, in many cases, be achieved by incorporating prior knowledge about the problem, captured in the form of constraints.

When modelling complex systems we are unavoidably faced with *imperfect* or *missing* information, especially in the measurement and information sciences. This may have several causes, but it is mainly due to

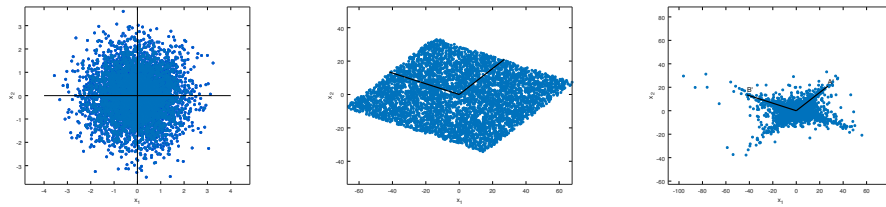


Figure 1: *Seeking structure in data.* Component analysis can be viewed as a family of data representation methods. The challenging task is to find *informative directions* in data space. These correspond to the column vectors of the observation (transformation, or ‘mixing’) matrix and form a new coordinate system. Their directions are non-orthogonal in general. (Left) Rotational invariance of the distribution of independent Gaussian random variables with equal variance. A scatterplot (point cloud) drawn from two such Gaussian sources illustrates the fact that there is not enough structure in the data in order to find characteristic directions in data space. Algebraically, we can only estimate the linear map up to an orthogonal transformation. (Center) Point cloud generated from a non-Gaussian distribution. (Right) The data cloud contains more structure in this case, which we want to exploit. In particular, the geometric shape of the point cloud of this figure is an example of a dataset that is *sparse* with respect to the coordinate axes shown by the two arrows.

- 30 • Lack of, or incompleteness in, our understanding or knowledge of the  
31 phenomena involved.
- 32 • The cost of obtaining and processing the vast amounts of information often  
33 needed for a more complete measurement of the phenomena.
- 34 • Inherent system complexity and stochasticity.

35 Probability theory is a conceptual and computational framework for reasoning  
36 under uncertainty. Probabilities model *uncertainty* regarding the occurrence of  
37 *random* events. Assigning probability measures on uncertain quantities reflects  
38 precisely our lack of information about the quantities at hand. According to  
39 Cox’s theorem [17], probability is the only consistent, universal logic framework  
40 for quantitatively reasoning under uncertainty. Moreover, probability theory  
41 offers a consistent framework for modelling and inference. Jaynes [38] viewed  
42 probability theory as a unifying tool for plausible reasoning in the presence  
43 of uncertainty. From a modeler’s point of view, the greatest practical advantage  
44 of probability theory is perhaps that it offers modularity and extensibility:  
45 probability theory acts as “glue” for linking different models together.

46 **2 Second order decompositions: Singular value**  
 47 **decomposition and Principal Component Anal-**  
 48 **ysis**

49 Singular value decomposition is an important method, originating in the Linear  
 50 Algebra and Numerical Analysis communities, with a vast repertoire of applica-  
 51 tions in the Applied Sciences and Data Analysis. It is often used as a subroutine  
 52 in more complicated models, and there exist versions of it that are very compu-  
 53 tationally efficient. We only present the basic ideas here; see [28] for a reference.

Let  $\mathbf{X}$  be a  $M \times N$  rectangular data matrix, where each row is a data point and each column is a “feature”,

$$\mathbf{X}_{M \times N} = \begin{bmatrix} x_{1,1} & \cdots & x_{1,N} \\ \vdots & \ddots & \vdots \\ x_{M,1} & \cdots & x_{M,N} \end{bmatrix},$$

and assume without loss of generality that  $M \geq N$ . The singular value decomposition (SVD) is a *factorization* of matrix  $\mathbf{X}$  such that

$$\mathbf{X} = \mathbf{U}\mathbf{S}\mathbf{V}^T, \quad (1)$$

54 where the  $M \times M$  orthogonal matrix  $\mathbf{U} = [\mathbf{u}_i]$  is called the left eigenvector  
 55 matrix of  $\mathbf{X}$ , and the  $N \times N$  orthogonal matrix  $\mathbf{V}^T = [\mathbf{v}_i^T]$  is its right eigenvector  
 56 matrix. The square roots of the  $N$  eigenvalues of the covariance matrices<sup>1</sup>  $\mathbf{X}\mathbf{X}^T$   
 57 and  $\mathbf{X}^T\mathbf{X}$  are the singular values of  $\mathbf{X}$ ,  $\sigma_i = \sqrt{\lambda_i}$ , forming the diagonal matrix  
 58  $\mathbf{S} = \text{diag}(\sigma_i)$ . The singular values are nonnegative and sorted in decreasing  
 59 order, such that  $\sigma_1 \geq \sigma_2 \geq \cdots \geq \sigma_N$ , forming the *spectrum* of  $\mathbf{X}$ .

The singular value decomposition of  $\mathbf{X}$  can be also written as

$$\mathbf{X} = \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{v}_i^T, \quad (2)$$

60 where  $\mathbf{u}_i$  is the  $i$ -th eigenvector of  $\mathbf{X}\mathbf{X}^T$  and  $\mathbf{v}_i$  is the  $i$ -th eigenvector of  $\mathbf{X}^T\mathbf{X}$ , as  
 61 above, and  $r \leq N$  is the rank of  $\mathbf{X}$ . In other words, a matrix,  $\mathbf{X}$ , can be written  
 62 as a linear superposition of its eigenimages, i.e. a sum of the outer products  
 63 of its left and right eigenvectors,  $\mathbf{u}_i \mathbf{v}_i^T$ , weighted by the square roots of the  
 64 eigenvalues,  $\sigma_i$ . The important fact here is that often relatively few eigenvalues  
 65 contain most of the ‘energy’ of matrix  $\mathbf{X}$ . Now if  $r < N$ , the energy of a data  
 66 matrix,  $\mathbf{X}$ , can be captured with fewer variables than  $N$ , since the relevant  
 67 information is contained in a lower-dimensional subspace of the measurement  
 68 space. This is a form of *dimensionality reduction*. Note that due to the presence  
 69 of noise in the data we may actually have  $r = N$ , though. In other words, in  
 70 practice all eigenvalues may be non-vanishing. This, however, also hints at a  
 71 *denoising* scheme in which one regards the smaller eigenvalues as corresponding

<sup>1</sup>Note that  $\mathbf{X}^T\mathbf{X} = \mathbf{V}\mathbf{S}^2\mathbf{V}^T$  and  $\mathbf{X}\mathbf{X}^T = \mathbf{U}\mathbf{S}^2\mathbf{U}^T$ .

72 to the noise and then forms a *truncated SVD*,  $\mathbf{X}_t = \mathbf{U}_t \mathbf{S}_t \mathbf{V}_t^T$ , where  $t < r$ .  
 73  $\mathbf{X}_t$  is the unique minimizer of  $\|\mathbf{X} - \mathbf{X}_t\|_F$  among all rank- $t$  matrices under the  
 74 Frobenius norm and also a minimizer (perhaps not unique) under the 2-norm.  
 75

76 **Remark 1** For many important applications, such as fMRI and other biosig-  
 77 nals, the signal of interest represents only a small part of what is measured  
 78 (Lazar, [43]; Calhoun et al., [10]), in terms of signal power. Consequently,  
 79 an optimization criterion that searches for components with maximum signal  
 80 power, such as PCA, will fail to recover the signals we are looking for. Methods  
 81 that exploit higher-order statistics in the data are therefore needed. Second-  
 82 order methods can still be very useful as a preprocessing step, however, e.g. for  
 83 dimensionality reduction, and are often used as such.

### 84 3 Higher-order decompositions: Independent Com- 85 ponent Analysis

86 In this section we review the independent component analysis (ICA) approach to  
 87 source separation, with an emphasis on the aspect of *non-gaussianity*. Method-  
 88 ological and review literature includes [57], [31], [58], [25]. Additional resources  
 89 are given below.  
 90

ICA is a family of data analysis methods that aims at decomposing datasets  
 into maximally statistically independent components. In the noiseless setting,  
 the observation model for linear ICA is

$$\mathbf{x} = \mathbf{A}\mathbf{s} \text{ ,} \quad (3)$$

where we have assumed that the observations have been de-meaned (i.e. we  
 have translated the coordinate system to the data centroid). ICA employs the  
 principle of *redundancy reduction* (Barlow, [5]) embodied in the requirement  
 of *statistical independence* among the components (Nadal and Parga, [50]).  
 In statistical language, this means that the joint density factorizes over latent  
 sources:

$$P(\mathbf{s}) = \prod_{l=1}^L P_l(s_l) \text{ ,} \quad (4)$$

91 where  $P(\mathbf{s})$  is the assumed distribution of the sources,  $\mathbf{s} = (s_1, \dots, s_L)$ , regarded  
 92 as stochastic variables, and  $p_l(s_l)$  are appropriate *non-Gaussian* priors. Non-  
 93 Gaussianity is the defining characteristic of the ICA family with respect to PCA.

94 We seek non-Gaussian sources for two, complementary, reasons:

- 95 • Identifiability,
- 96 • “Interestingness”.

97 Gaussians are not interesting since the superposition of independent sources  
 98 tends to be Gaussian. The concept of interestingness is directly exploited in the  
 99 related method of *Projection Pursuit* (Friedman and Tukey, [24]), where the  
 100 goal is to find the projection directions in a data set that show the least Gaus-  
 101 sian distributions. An important result relating to the probability densities  
 102 of the individual sources, due to Comon and based on the Darmois-Skitovitch  
 103 theorem<sup>2</sup> [18], formalizes the above and states that for analysis in independent  
 104 components at most one source may be Gaussian, in order for the model to  
 105 be estimable [16]. Geometrically, this indeterminacy of Gaussian point clouds  
 106 is due to the rotational invariance of the Gaussian distribution under orthog-  
 107 onal transformations (Hyvärinen, [34]). Gaussian point clouds are optimally  
 108 described in terms of the PCA decomposition method (Figure 1 (left) [Lewicki  
 109 and Sejnowski] (right)). This, geometric view of component analysis is a funda-  
 110 mental one in this paper.

111 A related concept is that of *linear structure* (Rao, [55]; Beckmann and Smith,  
 112 [7]). A vector,  $\mathbf{x}$ , is said to have a linear structure if it can be decomposed as  
 113  $\mathbf{x} = \boldsymbol{\mu} + \mathbf{A}\mathbf{s}$ , where  $\mathbf{s}$  is a vector of statistically independent random variables  
 114 and the matrix  $\mathbf{A}$  is of full column rank. Beckmann and Smith use results from  
 115 Rao [55] in order to ensure uniqueness of their ICA decomposition. In particular,  
 116 they use the fact that conditioned on knowing the number of sources and the  
 117 assumption of non-Gaussianity, there is no non-equivalent decomposition into a  
 118 pair  $(\mathbf{A}, \mathbf{s})$ , that is, there is no other decomposition with mixing matrix that is  
 119 not a rescaling and permutation of  $\mathbf{A}$ .

120 Equation (4) is equivalent to minimizing the *mutual information* among the  
 121 inferred sources<sup>3</sup> [8],

$$\begin{cases} \min I(s_1, \dots, s_L), & \text{where} \\ I(s_1, \dots, s_L) = \int p(\mathbf{s}) \log \frac{p(\mathbf{s})}{\prod_l p(s_l)} d\mathbf{s} \end{cases} \quad ,$$

or, equivalently, the “distance” between the distribution  $p(\mathbf{s})$  and the fully fac-  
 torized one,  $\prod_l p(s_l)$ , measured in terms of the Kullback-Leibler divergence,  
 $KL[p(\mathbf{s})||\prod_l p(s_l)]$ . This is defined as  $KL[p(x)||q(x)] = \mathbb{E}_{p(x)} \left[ \log \frac{p(x)}{q(x)} \right]$ . This

<sup>2</sup>The Darmois-Skitovitch theorem reads:

**Theorem 1 (Darmois-Skitovitch)** *Let  $\xi_1, \dots, \xi_n$  be independent random variables and let  $\alpha_i$  and  $\beta_i$ ,  $i = 1, \dots, n$  be nonzero real numbers such that the random variables  $\sum_{i=1}^n \alpha_i \xi_i$  and  $\sum_{i=1}^n \beta_i \xi_i$  are independent. Then the  $\xi_i$ 's are Gaussian.*

See, for example, V. Bogachev, ‘*Gaussian measures*’ [9], p. 13.

<sup>3</sup>For two stochastic variables  $X$  and  $Y$  to be independent, it is necessary and sufficient that their mutual information equals zero:

$$\begin{aligned} I(X, Y) &= H(X) + H(Y) - H(X, Y) \\ &= \int dX dY P_{X,Y}(X, Y) \log P_{X,Y}(X, Y) \\ &\quad - \int dX P_X(X) \log P_X(X) - \int dY P_Y(Y) \log P_Y(Y) = 0 \quad , \end{aligned}$$

where the quantity  $H(Z)$  is the ‘differential’ entropy of the random variable  $Z$ .

enables ICA algorithms to separate statistically independent sources, up to possible permutations and scalings of the components [16]. The mutual information (“redundancy”) can be equivalently computed as

$$I(s_1, \dots, s_L) = \left( \sum_{l=1}^L H(s_l) \right) - H(s_1, \dots, s_L) ,$$

122 where the first term at the RHS is the sum of the entropies of the individual  
 123 sources and the second the joint entropy of  $(s_1, \dots, s_L)$ . As shown by Bell  
 124 & Sejnowski (1995), independence can lead to separation because the method  
 125 exploits higher-order statistics in the data, something that cannot be done with  
 126 methods such as PCA.

In practice, many ICA algorithms minimize a variety of ‘proxy’ functionals. Bell and Sejnowski’s ICA approach uses the InfoMax principle (Linsker, [46]), maximizing *information transfer* in a network of nonlinear units (Bell & Sejnowski, [8]). Based on this, Bell and Sejnowski derive their very successful Infomax-ICA algorithm. The sources are estimated as

$$\hat{\mathbf{s}} = \mathbf{u} = \mathbf{W}\mathbf{x} , \quad (5)$$

where  $\mathbf{W}$  is the separating (unmixing) matrix that is iteratively learned by the rule

$$\mathbf{W} \leftarrow \mathbf{W} + \eta \left( \mathbf{I} - \mathbb{E}[\phi(\mathbf{u})\mathbf{u}^T] \right) \mathbf{W} , \quad (6)$$

127 where the vector valued map  $\phi(\mathbf{u}) = (\phi_1(u_1), \dots, \phi_L(u_L))$  is an appropriate  
 128 nonlinear function of the output,  $\mathbf{u}$ , such as a sigmoidal ‘squashing’ function,  
 129 applied component-wise. Popular choices are the logistic transfer function,  
 130  $\phi(u) = \frac{1}{1+e^{-u}}$ , and hyperbolic tangent,  $\phi(u) = \tanh(u)$ . The expectation oper-  
 131 ator,  $\mathbb{E}[\cdot]$ , is approximated by an average over samples in practice. Finally,  
 132 the factor  $\eta$  is an appropriate learning rate. The above equation incorporates  
 133 Amari et al.’s natural gradient descent approach [1]. Bell and Sejnowski show  
 134 that optimal information transfer, that is maximum mutual information be-  
 135 tween inputs and outputs, or equivalently maximum entropy for the output, is  
 136 obtained when highly-sloping parts of the transfer function are aligned with  
 137 high-density parts of the probability density function of the inputs.

Hyvärinen chooses to focus explicitly on non-Gaussianity and derives a fixed-point algorithm, dubbed FastICA [33]. Non-Gaussianity can be quantified using the *negentropy*,  $J$ ,

$$J(\mathbf{u}) = H(\mathbf{u}_{\text{Gauss}}) - H(\mathbf{u}) ,$$

where  $\mathbf{u}_{\text{Gauss}}$  is a Gaussian random variable with the same covariance as  $\mathbf{u}$ . The FastICA algorithm maximizes an approximation of  $J$  using the estimate

$$J(u_l) \approx \left\{ \mathbb{E}[G(u_l)] - \mathbb{E}[G(u_{\text{Gauss}})] \right\}^2 ,$$

where  $G(\cdot)$  is an appropriate nonlinearity, such as the non-quadratic function  $G(z) = z^4$ , and that is implicitly related to the source distributions (see below),

$u_{\text{Gauss}}$  is a standardized Gaussian r.v., and  $u_1, \dots, u_l, \dots, u_L$  are also of mean zero and unit variance. The unknown sources,  $\{u_l\}_{l=1}^L$ , are again estimated using the projections  $u_l = \mathbf{w}_l^T \mathbf{x}$ , where  $\mathbf{w}_l$  is the  $l$ -th separating vector (column of  $\mathbf{W}$ ), found by the iteration

$$\mathbf{w} \leftarrow \mathbb{E} [\mathbf{x}g(\mathbf{w}^T \mathbf{x})] - \mathbb{E} [g'(\mathbf{w}^T \mathbf{x})] \mathbf{w} ,$$

138 where  $g(\cdot)$  is the derivative of  $G(\cdot)$  and  $g'(\cdot)$  is the derivative of  $g(\cdot)$  and  $\mathbf{w}$  is  
 139 each time rescaled as  $\mathbf{w} \leftarrow \frac{\mathbf{w}}{\|\mathbf{w}\|}$ . For an application of the non-Gaussianity  
 140 principle to fMRI see the Probabilistic ICA algorithm of Beckman and Smith  
 141 [7].

### 142 ICA as Unfolding plus Rotation of a Dataset

An important result in the theory of ICA, with practical value, is that the ICA decomposition can be written as a factorization of an “*unfolding*” matrix times a *rotation* matrix. The former is usually implemented by pre-whitening (pre-sphering) the observations, such that  $\mathbb{E} [\tilde{\mathbf{x}}\tilde{\mathbf{x}}^T] = \mathbf{I}_D$ , where  $\tilde{\mathbf{x}}$  now denotes the whitened observations:

$$\tilde{\mathbf{x}} = \mathbf{W}_{\text{sph}} \mathbf{x} .$$

$\mathbf{W}_{\text{sph}}$  can be computed from the eigendecomposition of the data covariance matrix,  $\mathbf{C}_{\mathbf{x}\mathbf{x}} = \mathbb{E} [\mathbf{x}\mathbf{x}^T] \doteq \mathbf{U}\mathbf{\Lambda}\mathbf{U}^T$ , where the matrix  $\mathbf{U}$  is a unitary matrix<sup>4</sup> containing the eigenvectors of  $\mathbf{C}_{\mathbf{x}\mathbf{x}}$  and  $\mathbf{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_D)$  is the diagonal matrix of eigenvalues. Then the decomposition problem can be written (taking the “square root” and inverting) as

$$\tilde{\mathbf{x}} = \mathbf{\Lambda}^{-\frac{1}{2}} \mathbf{U}^T \mathbf{A} \mathbf{s} = \mathbf{W}_{\text{sph}} \mathbf{A} \mathbf{s} = \tilde{\mathbf{A}} \mathbf{s}, \quad \text{i.e.} \quad \mathbf{A} = \mathbf{W}_{\text{sph}}^{-1} \tilde{\mathbf{A}} .$$

That  $\mathbf{\Lambda}^{-\frac{1}{2}} \mathbf{U}^T$  spheres the data can be seen by simply performing the operations for  $\mathbb{E} [\tilde{\mathbf{x}}\tilde{\mathbf{x}}^T]$ , taking into account that  $\mathbf{U}$  is an orthogonal matrix [29]. The above whitening operation transforms the original data vectors to the space of the eigenvalues and rescales the axes by the singular values. Alternatively, one may use  $\mathbf{U}\mathbf{\Lambda}^{-\frac{1}{2}}\mathbf{U}^T$  for whitening, which maps the data back to the original space. This often makes further processing easier. In any case, since the whitening transformation removes any second-order statistics (correlations) in the data, learning the ICA matrix  $\tilde{\mathbf{A}}$  is equivalent to learning a pure orthogonal rotation matrix:

$$\mathbb{E} [\tilde{\mathbf{x}}\tilde{\mathbf{x}}^T] = \tilde{\mathbf{A}} \mathbb{E} [\mathbf{s}\mathbf{s}^T] \tilde{\mathbf{A}}^T = \tilde{\mathbf{A}} \tilde{\mathbf{A}}^T = \mathbf{I} .$$

### 143 3.1 Probabilistic Inference for ICA

Note that until now, while we have used probabilistic concepts to define information-theoretic quantities such as the negentropy and the mutual information, we have taken the view that the solution of the blind source separation problem can be

<sup>4</sup>If we restrict ourselves to the field of real numbers,  $\mathbb{R}$ , then the matrices  $\mathbf{U}$  become orthogonal matrices.

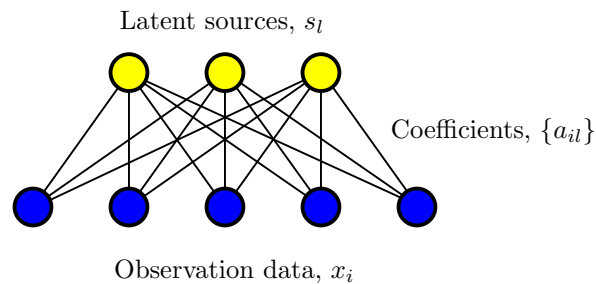


Figure 2: Graphical probabilistic model of the generative approach to component analysis. All models in this paper can be represented in this form.

achieved by transforming the observed signals through nonlinear functions in a bottom-up, filtering manner. Many classical component analysis algorithms, however, including ICA, can also be interpreted under the same probabilistic framework as top-down, *generative* models. This requires the construction of a density model. The model we consider here is the noisy transformation

$$\mathbf{s} \mapsto \mathbf{x} = \mathbf{A}\mathbf{s} + \underbrace{\boldsymbol{\varepsilon}}_{\text{noise}}, \quad (7)$$

where an  $L$ -dimensional vector of *latent variables*,  $\mathbf{s}$ , is linearly related to a  $D$ -dimensional vector of observations via the observation operator  $\mathbf{A}$ . Observation noise,  $\boldsymbol{\varepsilon}$ , may in general be added to the observations. In other words, the observed data is ‘explained’ by the unobserved latent variables, while the mismatch between the observations and the model predictions,  $\mathbf{x} - \mathbf{A}\hat{\mathbf{s}}$ , is explained by the additive noise. The fundamental equation of ICA, which we write again below,

$$P(\mathbf{s}) = \prod_{l=1}^L P_l(s_l), \quad (8)$$

144 can be seen as a modelling assumption, i.e. a *working hypothesis*, as a fac-  
 145 torization of a multi-dimensional distribution into a product of simpler one-  
 146 dimensional distributions, in another interpretation. Classical ICA models such  
 147 as Infomax ICA and FastICA assume noiseless and square mixing. This restric-  
 148 tion is removed in more recent algorithms. A representation of the generative  
 149 model for component analysis as a graphical probabilistic model is shown in  
 150 Fig. 2.

151 **Remark 2** *The generative model of Eqns (7), (8) defines a constrained prob-*  
 152 *ability distribution in data space. Referring back to Fig. 1, the “arms” of the*  
 153 *point-cloud are oriented along the directions of the “regressors”, which are en-*  
 154 *coded in the column vectors of the mixing matrix. Thus, when defining and*  
 155 *learning a probabilistic ICA model, we are in fact defining at least three*  
 156 *things: the source distributions, the mixing matrix, and the noise model, given*  
 157 *the constraints of Eqns 7 and 8.*



158 *This remark is important, as it gives an insight into why ICA algorithms are*  
 159 *so successful in decomposing certain types of data such as fMRI [19].*

In the general, noisy and non-square mixing case, one can formulate the penalized optimization problem (see e.g. [47], [11], [54], [61], and [59] for a nice concise review)

$$\hat{\mathbf{s}} = \underset{\mathbf{s}}{\operatorname{argmax}} \left\{ -\frac{1}{2\sigma^2} \|\mathbf{x} - \mathbf{A}\mathbf{s}\|^2 + \sum_{l=1}^L \log p_l(s_l) \right\}, \quad (9)$$

160 assuming spherical Gaussian noise,  $\boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_L)$ , for example, in order to  
 161 reconstruct the sources from the inputs at their most probable value.

As shown by MacKay [47] and Pearlmutter and Parra [54], Infomax-ICA can be interpreted as a maximum likelihood model. Assuming square mixing (i.e. as many latent dimensions as observations,  $L = D$ ), and invertibility of the mixing matrix, the separating matrix is  $\mathbf{W} = \mathbf{A}^{-1}$ . We can then immediately write down the probability of the data, as

$$p(\mathbf{x}) = |\det(\mathbf{J})| p(\mathbf{s}),$$

where  $\mathbf{J}$  is the Jacobian matrix of the transformation, with  $J_{li} = \frac{\partial s_l}{\partial x_i}$ . Under the linear model, and using the fundamental assumption of ICA, of mutual independence of the latent variables,  $p(\mathbf{s}) = \prod_{l=1}^L p(s_l)$ , we have

$$p(\mathbf{x}) = |\det(\mathbf{W})| \prod_l p(s_l).$$

Then, the log-likelihood of an *i.i.d.* data set,  $\mathcal{X} = \{\mathbf{x}_n\}_{n=1}^N$ , under the model can then be written as

$$\mathcal{L}(\boldsymbol{\theta}) \stackrel{\text{def}}{=} \log p(\mathcal{X}|\boldsymbol{\theta}) = \log \left( \prod_{n=1}^N p(\mathbf{x}_n|\boldsymbol{\theta}) \right) = N \log |\det(\mathbf{W})| + \sum_{n=1}^N \sum_{l=1}^L \log(p_l(\mathbf{w}_l^\top \mathbf{x}_n)),$$

162 where we have substituted  $s_{l,n}$  with  $u_{l,n} = \mathbf{w}_l^\top \mathbf{x}_n = \sum_{i=1}^D w_{l,i} x_{i,n}$ . The param-  
 163 eter vector,  $\boldsymbol{\theta}$ , here contains the matrix,  $\mathbf{A}$ , or equivalently the unmixing one,  
 164  $\mathbf{W} = \mathbf{A}^{-1}$ , since these are uniquely related in this case.

We can now derive a maximum likelihood algorithm for ICA via gradient descent, in order to learn the separating matrix,  $\mathbf{W}$ . Taking the derivative of  $\mathcal{L}(\boldsymbol{\theta})$  with respect to  $\mathbf{W}$  and using well-known derivative rules we finally find the learning rule

$$\frac{\partial}{\partial W_{li}} \mathcal{L}(\boldsymbol{\theta}) = A_{li} + z_l x_i,$$

165 where we have used the shorthand notation  $z_l = \phi_l(u_l)$ , where the ICA nonlin-  
 166 earity is the *score* function of the sources,  $\phi_l(s_l) = -\frac{\partial}{\partial s_l} \log p_l(s_l)$ , where  $p_l(s_l)$   
 167 are the *assumed* source priors. Multiplying with  $\mathbf{W}^\top \mathbf{W}$ , to make the algorithm  
 168 covariant [47], we get exactly the Infomax-ICA update rule, Eq. (6). Note that

169 the above multiplication is equivalent to using the ‘natural gradient’ approach of  
 170 Amari [1], a learning algorithm based on the concept of information geometry.  
 171 The FastICA algorithm can be also interpreted as an instance of the EM  
 172 algorithm [20], an iterative method for finding maximum likelihood or maxi-  
 173 mum a-posteriori solutions of statistical estimation problems. (See the “The  
 174 EM Algorithm” sidebar.) Lappalainen [41] derives it as an algorithm that fil-  
 175 ters Gaussian noise. This is an important interpretation, as it leads us to a  
 176 conceptually new framework for ICA, that of source separation via *denoising*.  
 177 Here, the term ‘denoising’ is interpreted as filtering out irrelevant information.  
 178 It is worth going through the main steps of the derivation.

#### The EM Algorithm

The general idea of the EM algorithm is to estimate the latent variables,  $\mathbf{Y}$ , and model parameters,  $\boldsymbol{\theta}$ , of a probabilistic model (which in this case are the sources,  $\mathbf{S}$ , and mixing matrix,  $\mathbf{A}$ , of the BSS problem, respectively), in two alternating steps. The ‘E’ (expectation) step computes the expectation of the log-likelihood with respect to the posterior distribution  $p(\mathbf{Y}|\mathbf{X}, \boldsymbol{\theta}^{(\tau)})$ , using the current ( $\tau$ th) estimate of the parameters,  $\boldsymbol{\theta}^{(\tau)}$ , giving the so-called ‘Q-function’,

$$Q(\boldsymbol{\theta}|\boldsymbol{\theta}^{(\tau)}) = \mathbb{E}_{\mathbf{Y}|\mathbf{X}, \boldsymbol{\theta}^{(\tau)}} [\log \mathcal{L}(\boldsymbol{\theta}; \mathbf{X}, \mathbf{Y})] ;$$

this is a function of  $\boldsymbol{\theta}$  only. (Recall that  $\mathbf{X}$  is observed and  $\boldsymbol{\theta}^{(\tau)}$  is temporarily fixed to its current point estimate.) The ‘M’ (maximization) step then computes the model parameters that maximize the expected log-likelihood,

$$\boldsymbol{\theta}^{(\tau+1)} = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} Q(\boldsymbol{\theta}|\boldsymbol{\theta}^{(\tau)}) .$$

This scheme is iterated until the algorithm converges. It can be shown that the EM algorithm is guaranteed to increase the observed data likelihood at each iteration [20].

179

Applying the above generic EM recipe, we can compute the maximum likelihood estimate of the mixing matrix of our ICA model as

$$\hat{\mathbf{A}} = (\mathbf{X} \mathbb{E}[\mathbf{S}]^T) (\mathbb{E}[\mathbf{S}\mathbf{S}^T])^{-1} ,$$

where the expected *sufficient statistics*<sup>5</sup> of the sources,  $\mathbb{E}[\mathbf{S}]$  and  $\mathbb{E}[\mathbf{S}\mathbf{S}^T]$ , are computed with respect to their posterior<sup>6</sup>. In the low sensor noise ( $\sigma^2 \rightarrow 0$ ) and square-mixing case of FastICA, Lappalainen approximates the posterior mean

<sup>5</sup>A sufficient statistic is the minimal statistic that provides sufficient information about a statistical model. Typically, the sufficient statistic is a simple function of the data, e.g. the sum of all the data points, sum of squares of the data points, etc.

<sup>6</sup>These are relationships that will become useful later in the paper as well.

of the sources as

$$\hat{\mathbf{s}} = \mathbb{E} [\mathbf{s} | \mathbf{A}, \mathbf{x}, \sigma^2 \mathbf{I}_D] \approx \mathbf{s}_0 + \sigma^2 (\mathbf{A}^\top \mathbf{A})^{-1} \boldsymbol{\phi}(\mathbf{s}_0) ,$$

180 where  $\mathbf{s}_0 \stackrel{\text{def}}{=} \mathbf{A}^{-1} \mathbf{x}$  and the function  $\boldsymbol{\phi}(\cdot)$  is defined as before, as the vector  
 181 of the logarithmic derivatives of  $p_l(\mathbf{s}_l)$ . For prewhitened data, this expression  
 182 simplifies even more, since  $\mathbf{A}$  is orthogonal, and therefore,  $(\mathbf{A}^\top \mathbf{A})^{-1} = \mathbf{I}_L$ .  
 183 Then  $\hat{\mathbf{A}} \approx \mathbf{A} + \sigma^2 \mathbf{X} \boldsymbol{\phi}(\mathbf{s}_0) / M$ .

Now Lappalainen makes the crucial observation that while the EM algorithm has not yet converged to the optimal values, the sources,  $\mathbf{s}_0$ , can be written as a “mixture”

$$\mathbf{s}_0 = \alpha \mathbf{s}_{\text{opt}} + \beta \mathbf{s}_G, \quad \text{with } \alpha^2 + \beta^2 = 1 ,$$

where the “noise”  $\mathbf{s}_G$  is mostly due to the other sources not having been perfectly unmixed. When far from the optimal solution, we have  $\beta \approx 1$  and  $\alpha \approx 0$ . Using an argument based on the central limit theorem, as the number of the other sources becomes large he then approximates the mixing matrix corresponding to those other sources as

$$\hat{\mathbf{a}}_G \approx \mathbf{a} + \sigma^2 \mathbf{X}_G \boldsymbol{\phi}(\mathbf{s}_{0G})^\top / L ,$$

where  $\mathbf{X}_G$  are Gaussian-distributed “sources” with the same covariance as  $\mathbf{X}$ , as is done in the standard FastICA algorithm, and the sources  $\mathbf{s}_{0G}$  are  $\mathbf{s}_{0G} \stackrel{\text{def}}{=} \mathbf{a}^\top \mathbf{X}_G$ . Then the update equation for the mixing matrix, normalized to unity, is estimated by

$$\hat{\mathbf{a}}_{\text{new}} = \frac{\hat{\mathbf{a}} - \hat{\mathbf{a}}_G}{\|\hat{\mathbf{a}} - \hat{\mathbf{a}}_G\|} \approx \frac{\sigma^2 [\mathbf{X} \boldsymbol{\phi}(\mathbf{s}_0)^\top - \mathbf{X}_G \boldsymbol{\phi}(\mathbf{s}_{0G})^\top] / L}{\|\hat{\mathbf{a}} - \hat{\mathbf{a}}_G\|} .$$

184 Lappalainen interprets the above E-step as *filtering* Gaussian noise.

185 The final step that will bring us to the standard FastICA is to note that the  
 186 term  $\mathbf{X}_G \boldsymbol{\phi}(\mathbf{s}_{0G})^\top / L$  is equal to  $\mathbf{a} \mathbf{s}_{0G} \boldsymbol{\phi}(\mathbf{s}_{0G})^\top / L$ , where the factor  $\mathbf{s}_{0G} \boldsymbol{\phi}(\mathbf{s}_{0G})^\top / L$   
 187 is constant, and therefore the numerator of the update equation becomes the  
 188 standard FastICA update,  $\hat{\mathbf{a}} - \hat{\mathbf{a}}_G = \mathbf{X} \boldsymbol{\phi}(\mathbf{s}_0)^\top - \mathbf{c} \mathbf{a}$ .

189 While Teh [59] computes the data likelihood in a maximum likelihood frame-  
 190 work, Knuth [39] uses a maximum a-posteriori framework. The latter allows us  
 191 to impose constraints on the model parameters as well. This was further ex-  
 192 plored in Hyvarinen and Karthikes in [32] in order to impose sparsity on the  
 193 *mixing* matrix.

194 Up to now we have either assumed equal number of sources and sensors  
 195 or we have implicitly assumed that their number is somehow given. Roberts  
 196 [56] derives a Bayesian algorithm for ICA under the evidence framework that  
 197 estimates the most probable number of sources as a model order estimation  
 198 problem. The evidence framework, as applied in [56], makes a local Gaussian  
 199 approximation to the likelihood conditioned on the mixing matrix using a nested  
 200 Laplace approximation, but takes into account the local curvature by estimating  
 201 the Hessian. Due to computational reasons, this is approximated by a diagonal

202 matrix here, setting the off-diagonal elements to zero. The noise width, regarded  
203 as a *hyperparameter*, is computed at its maximum likelihood value.

204 Finally, Choudrey et al. [15] and Miskin and MacKay [49] propose a fully  
205 Bayesian approach to ICA using a variational, *ensemble learning* approach under  
206 a mean-field approximation. They use a flexible source model based on mixtures  
207 of Gaussians and perform model order estimation using a variety of techniques.

208 We can now select an appropriate functional form for the individual marginal  
209 distributions,  $p_l(s_{l,n})$ , based on our prior knowledge about the problem, as was  
210 done in the original formulation of InfoMax ICA of Bell & Sejnowski for the  
211 separation of speech signals, for example. The source model should model  
212 the real source distributions as accurately as possible. Many natural signals  
213 exhibit characteristic amplitude distributions, which can provide some guidance  
214 and indeed should be exploited when possible. This allows us to utilize fixed  
215 source models in our separation algorithms. Bell and Sejnowski, for example,  
216 use several nonlinearities (recall that these are uniquely related to the assumed  
217 PDFs of the sources), such as  $1/(1+e^{-u_i})$ ,  $\tanh(u_i)$ ,  $e^{-u_i^2}$ , etc., as well as propose  
218 general-purpose ‘score functions’ (see Figure 2 of Ref. [8]) in their Infomax-ICA  
219 algorithm. FastICA uses nonlinearities such as  $u_i^3$ ,  $\tanh(\alpha u)$ ,  $u_i e^{-\alpha u_i^2/2}$ , and  
220  $u_i^2$ . However, this is not always possible. The problems that can arise from an  
221 incorrect latent signal model and possible solutions are discussed in section 4.

## 222 4 The Importance of using Appropriate Latent 223 Signal Models

224 Many classical ICA algorithms, such as Infomax-ICA and FastICA, allow the  
225 plug-in setting of the respective nonlinearity function in the system, as men-  
226 tioned above. For successful separation, the form of the nonlinearity must  
227 somehow *match*, as far as possible, the underlying (unknown) statistical prop-  
228 erties of the sources, such as their super- or sub-gaussianity. This was first stated  
229 as “matching the neuron’s input-output function to the expected distribution  
230 of the signals” in [8]. Since the estimating equations for the mixing matrix and  
231 sources are coupled, the functional form of the nonlinearity is critical for their  
232 correct estimation: an incorrect choice of nonlinearity will lead to an incorrect  
233 estimation of the (un-)mixing matrix, which will map the observations back  
234 to the source space incorrectly, etc. Cardoso [12] gives a compelling example  
235 of how estimation can go wrong. Another example of how classical ICA fails  
236 in separating sources in an image processing context is given in Fig. 4 (from  
237 Tonazzini et al., [60]).

238 **Remark 3** *Tonazzini et al. use a Markov random field in order to impose an*  
239 *image prior. However, the images of Fig. 4 (left) are actually also prime exam-*  
240 *ples of sparse sources. In [27] and [19], an extensive study of how justified and*  
241 *robust are ICA algorithms for functional MR imaging of the brain was conducted*  
242 *and various simulations of fMRI “brain” activations under well-controllable sit-*  
243 *uations with shapes similar to that of ref. [60] were performed that highlighted*

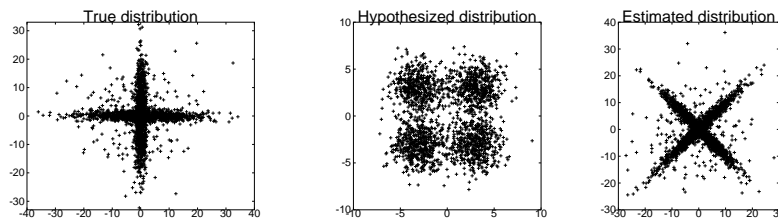


Figure 3: Effect of an incorrect source model specification [12]. Left: true distribution; Middle: Hypothesized distribution; Right: Estimated distribution.

244 *the need for alternative decomposition algorithms that are effective for fMRI,*  
 245 *based on sparsity.*

It can be shown that the Infomax-ICA as well as the FastICA algorithms are instances of maximum likelihood estimation [47], [54], [30], [41]. Under this interpretation, one can see that the nonlinearity,  $\phi(\cdot)$ , is actually the logarithmic derivative of the (hypothesized) probability density of the sources (the ‘score’ function): for the  $l$ -th source,  $\mathbf{s}_l$ ,

$$l : \quad \phi_l([\mathbf{W}\mathbf{x}]_l) = -\frac{\partial}{\partial \mathbf{s}_l} \log p_l(\mathbf{s}_l) = -\frac{p'_l(\mathbf{s}_l)}{p_l(\mathbf{s}_l)},$$

246 where the symbol  $\mathbf{W}$  denotes the separating operator from observation space  
 247 to source space and  $\mathbf{x}$  is an observation. In other words, in a perfect match the  
 248 nonlinearity is exactly the cumulative distribution function of the sources. Of  
 249 course we do not know the *actual* source PDFs, since the sources themselves  
 250 are unobserved, but we may try to estimate them from the data. For this  
 251 purpose, we can employ a parameterized model source PDF,  $p_l(\mathbf{s}_l; \theta_{\mathbf{s}_l})$ , and  
 252 *learn*, instead of fix, its parameters,  $\theta_{\mathbf{s}_l}$ , from the data. A flexible prior that is at  
 253 the same time mathematically tractable is a mixture distribution. Lawrence and  
 254 Bishop [42] uses a Mixture of Gaussians (MoG) prior for ICA, albeit in a fixed  
 255 form. Attias [3] has used MoGs as source models for blind source separation  
 256 under a maximum likelihood framework, leading to a flexible algorithm dubbed  
 257 ‘Independent Factor Analysis’ (IFA). Choudrey et al. [15] and Lappalainen [40]  
 258 use the same prior under a Bayesian ensemble learning approach, i.e. with a  
 259 factorized posterior (the so-called ‘naive’ mean-field method).

## 260 5 Sparse Decompositions

261 As noted by Cardoso [13], non-Gaussianity is not the only possible route to inde-  
 262 pendent component analysis, and indeed to blind source separation in general;  
 263 other possibilities also exist—including exploiting non-stationarity and time-  
 264 correlation in signals. Such a different paradigm, *sparsity*, in combination with  
 265 doing away with the assumption of independence, will be explored next.

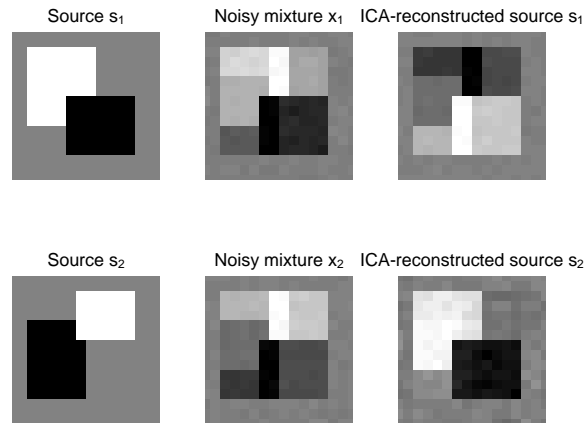


Figure 4: Effect of an incorrect source model specification for a blind image separation problem. Left: true sources; Middle: noisy mixtures; Right: Estimated sources from ICA. The model clearly fails to recover the sources. In particular, one of the sources is not recovered at all.

## 266 5.1 Parsimonious representation of data

267 ICA works well for a variety of blind source separation problems. However, in  
 268 order for the decomposition to make sense the true sources must themselves  
 269 indeed be (nearly) independent. This may make sense in the separation of voice  
 270 signals that are independently generated by people with no interaction among  
 271 them, for example. For other problems, however, searching for components  
 272 that are maximally independent may not be so meaningful. Recently, another  
 273 paradigm for BSS, and inverse problems in general, *sparsity*, has emerged as an  
 274 alternative. Sparsity refers to the property of a representation to form compact  
 275 encodings of signals, data, or functions, using a small number of basis functions.  
 276 Those basis functions are used as “building blocks” to build more complex  
 277 signals.

278 There has been a variety of algorithms for sparse representation, or sparse  
 279 coding, originating from the computational neuroscience and neural networks  
 280 communities as well as several others from a signal processing perspective.  
 281 Sparse decomposition, and ways to impose sparsity constraints, has recently  
 282 also been a topic of much research in the statistics and machine learning liter-  
 283 ature.

284

285 **Sparse coding.** In the study of the visual system, Field [23] proposed sparsity  
 286 as an organization principle of the visual receptive field. He conjectured that

287 populations of neurons optimize the representation of their visual environment  
 288 by forming *sparse representations* of natural scenes, a hypothesis that has high  
 289 biological plausibility since it is based on the general idea of a system using  
 290 its available resources efficiently. According to his theory, the visual system  
 291 performs efficient coding of natural scenes in terms of natural scene statistics  
 292 by finding the sparse structure available in the input. Field's theory directly  
 293 reflects the principle of *redundancy reduction* of Barlow [5], [6].

**Dictionary learning.** Olshausen and Field [51] further test the above theory, seeking experimental evidence for sparsity in the primary visual cortex (V1) by building a predictive (mathematical) model of sparse coding. In their model, images are formed as a linear combination of local basis functions with corresponding activations that are as sparse as possible. These bases model the V1 receptive fields and form overcomplete sets *adapted* to the statistics of natural images. Olshausen and Field's model is an early example of *dictionary learning*. Formally, the model of Olshausen and Field is described by:

$$\mathbf{x}_p \simeq \sum_i a_{p,i} \phi_i ,$$

where  $\mathbf{x}_p$  is an image "patch" (i.e. a small image window) and  $\{\phi_i\}$  are the underlying basis elements. A network representation of their model, **Sparsenet**, is shown in Fig. 5. They proposed the following objective:

$$\mathcal{I}(\Phi) = \min_{a_{p,i}} \left\{ \sum_p \left\| \mathbf{x}_p - \sum_i a_{p,i} \phi_i \right\|^2 + \lambda \sum_i \log(1 + a_{p,i}^2) \right\} ,$$

294 to be minimized over bases,  $\Phi$ , learned by searching for bases that optimized the  
 295 sparsity of the coefficients,  $\{a_{p,i}\}$ , (subject to appropriate scale normalization  
 296 of  $\{\phi_i\}$ ). In general, the basis set can be overcomplete. That is, the number  
 297 of bases,  $|\Phi|$ , can be greater than the dimensionality of the 'input' data space,  
 298  $D$  (see for example [52]). The reason for this is that the 'code' can be more  
 299 sparse if one allows an overcomplete basis set, as the algorithm can select the  
 300 bases that better match the structures contained in the signal (the "active"  
 301 elements). See also Asari, [2]. As shown in Fig. 6 this objective results in highly  
 302 sparse distributions for the coefficients. Astonishingly, the learned receptive  
 303 fields (filters) have properties that resemble the properties of natural simple-cell  
 304 receptive fields, that is they are spatially localized, oriented and bandpass, i.e.  
 305 selective to structure at different spatial scales (Fig. 7).

306 In the signal processing community, Mallat and Zhang [48] proposed a  
 307 greedy algorithm analogous to the projection pursuit in statistics, called 'matching  
 308 pursuit', that iteratively finds the best matching projections of signals onto  
 309 a fixed overcomplete dictionary of time-frequency 'atoms'. Linear combinations  
 310 of those atoms form compact representations of the given signal.

311

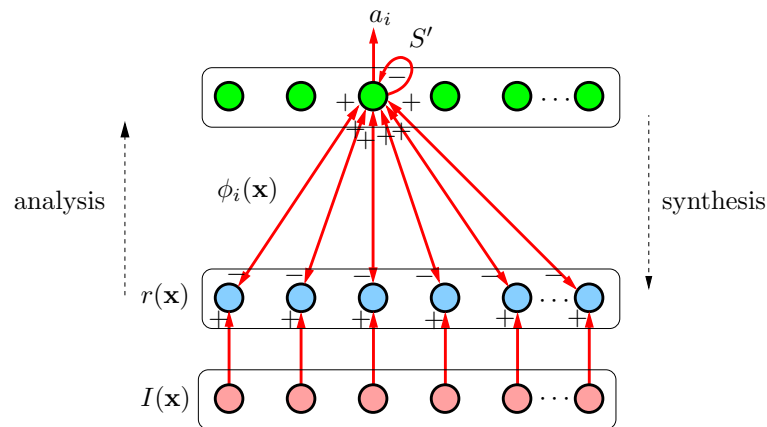


Figure 5: The Olshausen and Field model [51] as a neural network, **Sparsenet**. The inputs to the network are images,  $I(\mathbf{x})$ , where  $\mathbf{x}$  denotes picture elements (pixels) over an image domain,  $\Omega$ , and the outputs are the coefficients of the representation,  $a_i$ . The symbol  $r(\mathbf{x})$  is the residual image,  $r(\mathbf{x}) = I(\mathbf{x}) - \sum_i a_i \phi_i(\mathbf{x})$ . Each output neuron evolves according to the differential equation  $\dot{a}_i = \sum_{\mathbf{x} \in \Omega} \phi_i(\mathbf{x}) r(\mathbf{x}) - \lambda S'(a_i)$ , where the derivative of the sparsity activation function  $S(\cdot)$  induces non-linear self-inhibition, and the multiplier  $\lambda \geq 0$  is a regularization parameter. This enforces *sparsity*, as it drives activities towards zero. The regularization parameter balances the first, data fidelity term, which ensures accurate reconstruction. During the ‘analysis’ (“filtering”) phase, a given image,  $I(\mathbf{x})$ , is decomposed in a dictionary,  $\Phi$ , and its corresponding coefficients,  $a_i$ , are computed. During the ‘synthesis’ phase a learned dictionary predicts an estimate of an image,  $\hat{I}(\mathbf{x})$ , with residuals  $\mathbf{r}(\mathbf{x})$ . The optimal value of each  $a_i$  is determined from the corresponding equilibrium solution.



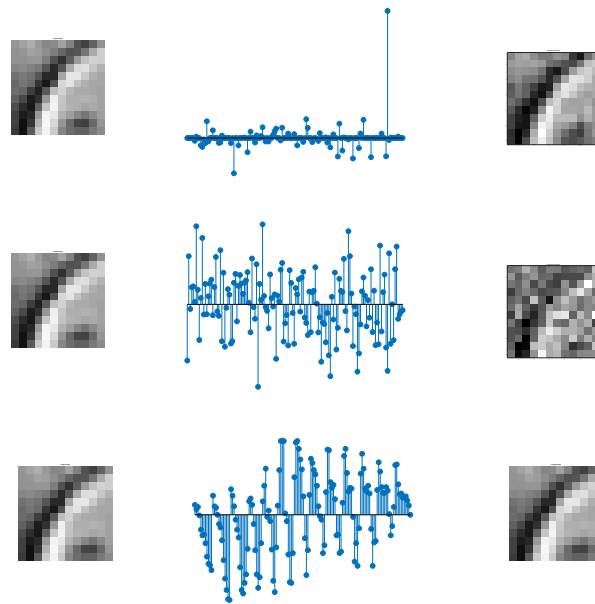


Figure 6: Activities,  $a_i$ , resulting from the model of Olshausen and Field [51]. The input image on the left is reconstructed from learned bases using their algorithm. Note how the coefficients  $a_i$  resulting from the model (first row) are highly sparse, compared to reconstructing the image patch using random bases (second row) or pixel (canonical) bases (third row). The canonical basis offers no compression at all, as it is merely a copy of the original image.

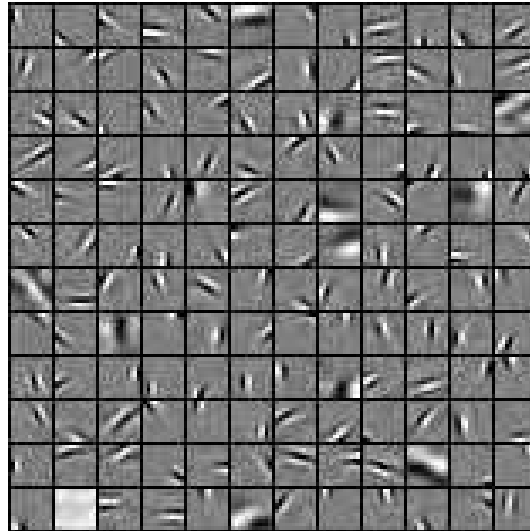


Figure 7: Learned receptive fields (filters) from the sparse coding algorithm of Olshausen and Field **Sparsenet** [51]. These filters exhibit properties of simple-cell receptive fields such as locality, orientation and spatial selectivity.

312 **Geometric interpretation of sparse representation.** A geometric inter-  
 313 pretation of sparse representation is depicted in Fig. 8. Each data vector can be  
 314 viewed as a point in a  $D$ -dimensional vector space, the whole dataset forming  
 315 a cloud of points. We now seek a linear transformation of the dataset such  
 316 that the inferred “projections” on to the new coordinate system defined by the  
 317 column vectors of the learned transformation matrix,  $\mathbf{A} = [\mathbf{a}_l]_{l=1}^L$ , are as sparse  
 318 as possible.

319 Note that it is the sparseness of the components (and the selection of a  
 320 suitable model prior) that drives learning of the new representation (unmix-  
 321 ing) directions. This sparseness is reflected in the *shape* of the point-cloud:  
 322 referring to the above figure (where  $D = L \doteq 2$ ), sparse data mapped in to  
 323 the latent space produce a highly-peaked and heavy-tailed distribution for both  
 324 axes (Fig. 8 (lower right)). This is indeed a result of the sparseness property of  
 325 the dataset: the two ‘arms’ of the sparse data cloud are *tightly packed* around  
 326 the directions of the unmixing vectors,  $\mathbf{a}_l$ . Algebraically, this means that for  
 327 a particular point,  $n$ , either the coefficient  $s_{1,n}$  ( $l = 1$ ) or the coefficient  $s_{2,n}$   
 328 ( $l = 2$ ) is almost zero, as the particular datum is well described by the  $\mathbf{a}_2$  or  
 329 the  $\mathbf{a}_1$  “regressor”, respectively. On the contrary, non-sparse data will typically  
 330 produce a projection that corresponds to a “fat” empirical histogram, as shown  
 331 in Fig. 8 (upper-right).

<sup>7</sup>Field studied the statistics of natural scenes and their relation to computer vision and perception in [23]. The ‘state-space’ in this context is a state-space of neural activation amplitudes.

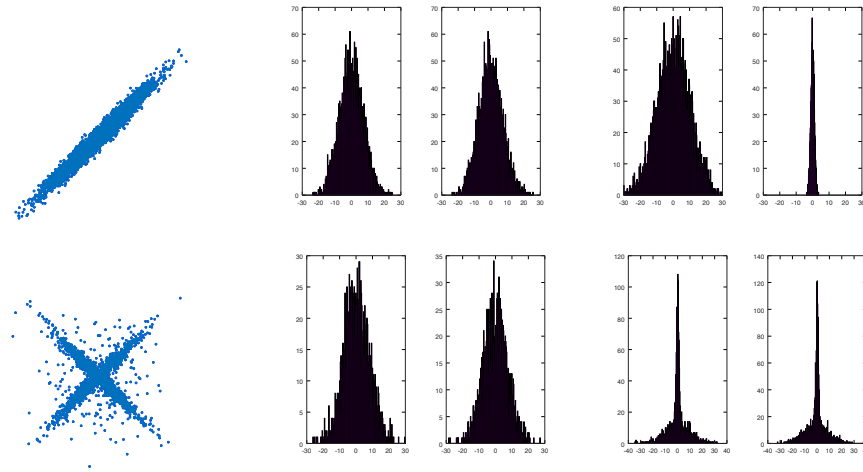


Figure 8: *Geometric interpretation of sparse representation.* State-spaces (in the terminology of Field<sup>7</sup>[23]) and projections of two datasets, one sparse (lower row) and the other non-sparse (upper row), are shown. Each dataset, plotted in the measurement coordinate system,  $xy$ , produces a point cloud (left part of the figure) — for visualization purposes, both observation and latent dimensionalities are equal to  $D = L = 2$  in this figure. By projecting the point clouds on to each coordinate we can produce the corresponding empirical histograms of ‘state’ amplitudes (middle part of the figure). We now seek a linear transformation to a latent space,  $uv$ , such that it optimizes some suitable criterion (this is shown in the right part of the figure). Sparse data mapped in the latent space produce heavy-tailed distributions for both latent dimensions (lower right), while for non-sparse data this is not the case (upper right).

With respect to the soft clustering view of component analysis (Miskin, [36]), discussed in the Introduction of the paper, if the data vectors are sufficiently sparse, their images on the unit hypersphere  $\mathbb{S}^{D-1}$ , i.e. the radial sections of their position vectors with the unit hypersphere, mapped as

$$\mathbf{x}_n \in \mathbb{E}^D \mapsto \hat{\mathbf{x}}_n \in \mathbb{S}^{D-1},$$

332 where the projection operator  $P : \mathbf{u} \mapsto \hat{\mathbf{u}} = \frac{\mathbf{u}}{\|\mathbf{u}\|}$  maps vectors along their radii,  
 333 *concentrate* around the unit vectors  $\{\hat{\mathbf{a}}_l\}_{l=1}^L$ ; see Fig. 9 and Ref. [62]. While  
 334 Miskin did not use this property per se for sparse decomposition, one can design  
 335 separation algorithms that exploit it [45].

## 336 5.2 Sparse Decomposition of Data Matrices

337 Inspired by the model of Olshausen and Field, Donoho [21] first points out  
 338 the connection and differences between the two lines of research, independent  
 339 component analysis and sparse decompositions, and he promotes the idea of

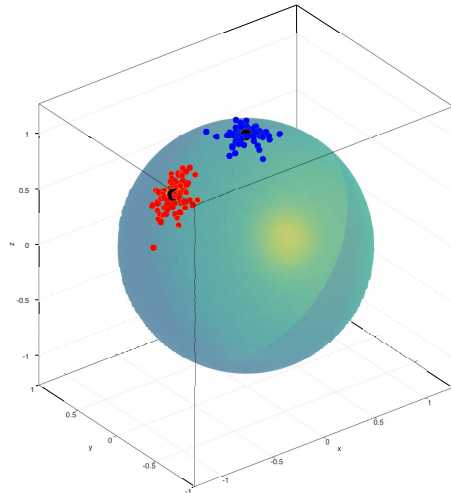


Figure 9: *Clustering of a sparse set of points on the unit hypersphere,  $\mathbb{S}^{D-1}$ , embedded in a  $D$ -dimensional space. The points cluster around the direction vectors corresponding to the columns of the mixing matrix.*

340 sparsity, overcompleteness, and optimal atomic decompositions as a better goal  
 341 than independence. He provides a rationale of why sparsity is a more plau-  
 342 sible principle, being “intrinsically important and fundamental”, due to both  
 343 biological and modelling reasons. Regarding the former, he too cites the ex-  
 344 tremely efficient sparse representation achieved by the human visual system,  
 345 and its higher compression performance compared to the best engineered sys-  
 346 tems. With respect to the latter, he notes that independence is inherently a  
 347 probabilistic assumption and of unknown interpretability (with respect to vi-  
 348 sion) because natural images are composed by occlusion. Occlusion inevitably  
 349 creates *dependent* components. He finally suggests that one of the future chal-  
 350 lenges of ‘sparse components analysis’ would be to search over spaces of objects  
 351 of much larger scale than the image patches of Olshausen and Field.

352 It turns out (see Olshausen, [52]) that the Infomax-ICA algorithm becomes,  
 353 in fact, a special case of the sparse linear algorithm of Olshausen and Field  
 354 when there is an equal number of basis functions/latent dimensions and inputs,  
 355 the  $\phi_i$ s are linearly independent, and there is no observation noise. In this case,  
 356 there is a unique set of coefficients  $\{a_i\}$  that is the root of  $\|\mathbf{X} - \Phi\mathbf{a}\|$ , and we  
 357 can write  $\mathbf{a}$  as  $\mathbf{a} = \mathbf{W}\mathbf{X}$ , where  $\mathbf{W} = \Phi^{-1}$  (note that based upon the above  
 358 assumptions,  $\Phi$  becomes invertible). If, in addition, the ICA nonlinearity is  
 359 chosen to be the cumulative density function of the sparse components, then

360 the sparse algorithm gives exactly the algorithm of Bell and Sejonwski. The  
 361 point here is actually to show that sparsity constraints can lead to separation.  
 362 Many researchers have indeed shown that this can be indeed the case. Indeed,  
 363 as pointed out by Li, Cichocki and Amari [45],

364 **Remark 4** *Sparse decompositions of data matrices can be used for the blind*  
 365 *source separation problem.*

They provide various examples from simulations and EEG data analysis that demonstrate the performance of sparse decompositions in signal separation. Li, Cichocki and Amari performed a sophisticated mathematical analysis for the case of sparse representation of data matrices under the  $\ell_1$  prior, for given basis matrices. They tackle the two-step decomposition problem of learning the base matrix first, via clustering, and then estimating the coefficients of the decomposition. If  $\mathbf{X}$  is a data matrix and  $\mathbf{A} = \{\mathbf{a}_l\}$  is a given basis, Li et al. start from the mathematical model shown below:

$$\min \left\{ \underbrace{\sum_{l=1}^L \sum_{n=1}^N |s_{ln}|}_{S(\mathbf{s})} \mid \text{subject to } \mathbf{A}\mathbf{s} = \mathbf{X} \right\}, \quad (10)$$

366 with  $S(\cdot)$  the *sparsity function* on the sources. This particular case of optimiza-  
 367 tion problem can then be solved using linear programming. While the  $\ell_0$ -norm  
 368 solution is the sparsest one in general, its optimization is a non-trivial combina-  
 369 torial problem. Li et al. show that, for sufficiently sparse signals, the solutions  
 370 to the problem of sparse representation of data matrices that are obtained using  
 371 the  $\ell_0$  and  $\ell_1$  norms are equivalent. This fact was previously shown by Donoho  
 372 and Elad [22] but Li et al. [45] give a less strict sparseness ratio (i.e. the ratio  
 373 of zero versus non-zero elements).

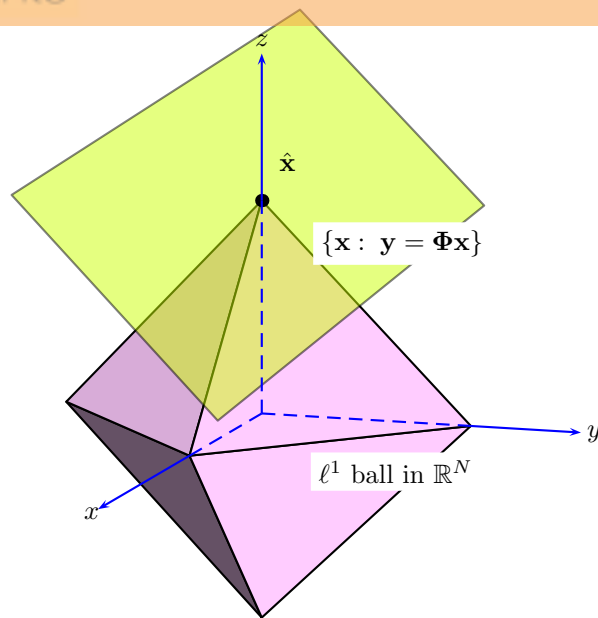
374

**Uniqueness.** Importantly, Li et al. [45] also show that the above problem has a *unique* solution. While in general there would be an infinite number of solutions for the underdetermined system of equations

$$\mathbf{A}\mathbf{s} = \mathbf{x},$$

where the  $D \times L$  matrix  $\mathbf{A}$  (observation operator) with  $L > D$  maps the unknown signal  $\mathbf{s}$  in to the observed signal  $\mathbf{x}$ , the sparsity constraint makes the particular linear inverse problem well-posed. A geometric interpretation of why  $\ell_1$ -type sparsity regularization works well for signal recovery under sparsity constraints is shown in Fig. 10. We want to find the optimal  $\mathbf{x}$  as the minimum-norm vector that satisfies the constraint  $\mathbf{x} = \mathbf{A}\mathbf{s}$ , i.e. such that the hyperplane does not intersect the  $\ell_1$  ball. More generally, the problem can be stated (in the deterministic framework) as:

$$\min_{\mathbf{s}} \left\{ \|\mathbf{s}\|_1 : \|\mathbf{A}\mathbf{s} - \mathbf{x}\| < c \right\}$$



1

Figure 10: *Why  $\ell_1$  works: A geometric intuition into sparse priors.* We seek the sparsest vector  $\mathbf{x} \in \mathbb{R}^N$  under the  $\ell_1$  norm, in this case, that satisfies the linear constraint  $\mathbf{y} = \Phi \mathbf{x}$ , where  $\Phi$  is a dictionary. The  $\ell_1$  penalty corresponds geometrically to a cross-polytope (the ' $\ell_1$  ball' in  $\mathbb{R}^N$ ) and the linear constraint to a hyperplane. The shape of the polytope dictates the form of the solution. The optimal vector,  $\hat{\mathbf{x}}$ , is the one that touches the hyperplane without the latter intersecting the cross-polytope. Mathematically, this is the solution to the problem  $\hat{\mathbf{x}} = \operatorname{argmin}_{\mathbf{y}=\Phi \mathbf{x}} \|\mathbf{x}\|_1$ . As can be seen from the figure, the inclusion of  $\ell_1$  norm necessarily drives all components of  $\mathbf{x}$  but one towards zero, leading to sparse solutions.

375 (Chen and Haykin, [14]), where  $\mathbf{x}$  can be a “corrupted” (noisy, blurred, etc)  
 376 version of the original signal and  $c$  is a positive scalar constant that plays a  
 377 role similar to the noise variance in the probabilistic framework (Li et al., [45]).  
 378 In this case, the hyperplane becomes an orthotope (hyperrectangle), defining a  
 379 “zone” in which the vertex of the  $\ell_1$  ball must fall. In addition, Li et al. [45]  
 380 use  $k$ -means clustering to get an estimate of the basis, which is then used  
 381 in a linear programming algorithm in order to estimate the coefficients of the  
 382 representation.

### 383 5.2.1 Probabilistic Solutions

Lewicki and Sejnowski [44], introduce a probabilistic method for sparse over-  
 complete representations. A Laplacian prior on the coefficients of the basis was  
 used,  $p(s_l) \propto e^{-\theta|s_l|}$ , enforcing parsimonious representations. They then pro-  
 pose a gradient optimization scheme for maximum a-posteriori (MAP) learning.  
 For the linear model  $\mathbf{x} = \mathbf{A}\mathbf{s} + \boldsymbol{\varepsilon}$ , with Gaussian observation noise with variance  
 $\sigma^2$ , we seek the most probable decomposition coefficients,  $\hat{\mathbf{s}}$ , such that

$$\hat{\mathbf{s}} = \underset{\mathbf{s}}{\operatorname{argmax}} \left\{ p(\mathbf{x}|\mathbf{A}, \mathbf{s})p(\mathbf{s}) \right\} . \quad (11)$$

The probability of a single data point is obtained by integrating out the unknown  
 signals,  $\mathbf{s}$ :

$$p(\mathbf{x}|\mathbf{A}) = \int p(\mathbf{x}|\mathbf{A}, \mathbf{s})p(\mathbf{s})d\mathbf{s} .$$

In order to derive a tractable algorithm, they make a Laplace approximation to  
 the data likelihood, by assuming that the posterior is Gaussian around the poste-  
 rior mode. This involves computing the Hessian  $\mathbf{H} = \nabla_{\mathbf{s}}\nabla_{\mathbf{s}} \{-\log [p(\mathbf{s})p(\mathbf{x}|\mathbf{A}, \mathbf{s})]\} =$   
 $\frac{1}{\sigma^2}\mathbf{A}^T\mathbf{A} - \nabla_{\mathbf{s}}\nabla_{\mathbf{s}} \log p(\mathbf{s})$ . To make a smooth approximation of the derivative of  
 the log-prior, and a diagonal approximation to the Hessian, they then take  
 $p(s_l) \approx \cosh^{-\theta/\beta}(\beta s_l)$ , which asymptotically approximates the Laplacian prior  
 for  $\beta \rightarrow \infty$ . Moreover, a low noise level is assumed. The above approximations  
 finally lead to the gradient learning rule

$$\Delta\mathbf{A} = \mathbf{A}^T\mathbf{A} \nabla_{\mathbf{A}} \log p(\mathbf{x}|\mathbf{A}) \approx -\mathbf{A} (\mathbf{I} + \mathbf{z}\hat{\mathbf{s}}^T) ,$$

384 where, again,  $z_l = \partial \log p(s_l)/\partial s_l$ . Note that this has the same functional form  
 385 as the Infomax-ICA learning rule, however the basis matrix is generally non-  
 386 square in this case. In contrast to the standard ICA learning rule, and where  
 387 the sources are estimated simply by  $\mathbf{s} = \mathbf{W}\mathbf{x}$ , where the unmixing matrix is  
 388  $\mathbf{W} = \mathbf{A}^{-1}$ , here we must use a nonlinear optimization algorithm in order to  
 389 estimate the coefficients, using Eq. (11). Due to the low-noise assumption, the  
 390 level of the observation noise is not estimated from the data and has to be set  
 391 manually. Lewicki and Sejnowski’s algorithm, however, is faster in obtaining  
 392 good approximate solutions than the linear programming method and is more  
 393 easily generalizable to other priors.

Girolami [26] proposes a variational method for learning sparse representations. In particular, his method offers a solution to the problem of analytically integrating the data likelihood, for a range of heavy-tailed distributions. Starting from the heavy-tailed distribution  $p(s) \propto \cosh^{-\frac{1}{\beta}}(\beta s)$ , he derives a variational approximation to the Laplacian prior by introducing a variational parameter,  $\xi = (\xi_1, \dots, \xi_L)$ , such that the prior  $p(\mathbf{s}) = \prod_{l=1}^L \exp(-|s_l|)$  becomes  $p(\mathbf{s}; \xi)$ , with  $\mathbf{s}|\xi \sim \mathcal{N}(\mathbf{s}; \mathbf{0}, \mathbf{\Lambda})$  and  $\mathbf{\Lambda} = \text{diag}(|\xi_l|)$ . Then  $p(\mathbf{s})$  is the supremum

$$p(\mathbf{s}) = \sup_{\xi} \left\{ \left[ \prod_{l=1}^L \varphi(\xi_l) \right] \mathcal{N}(\mathbf{s}; \mathbf{0}, \mathbf{\Lambda}) \right\},$$

394 with  $\varphi(\xi) \rightarrow \exp(-\frac{1}{2}|\xi|)\sqrt{2\pi|\xi|}$  as  $\beta \rightarrow \infty$ . The above is derived using a variational argument and using convex duality [37], [53]. In essence, what this  
 395 approximation means is that, at each point of its domain, the intractable prior  
 396 is lower-bounded tightly by a best-matching Gaussian with width parameter  $\xi$ ,  
 397 with this variational parameter being estimated by the algorithm along with the  
 398 model parameters. Using the above, the posterior takes a Gaussian form. This  
 399 enables him to derive an EM algorithm in order to infer the sparse coefficients  
 400 and learn the overcomplete basis vectors of the representation. Girolami applies  
 401 his sparse representation algorithm to the problem of overcomplete source separation  
 402 and achieves superior results compared to the algorithm of Lewicki and Sejnowski.  
 403  
 404

405 The problem of sparsely representing a data matrix described above is  
 406 a special case of the more general problem of recovering latent signals that  
 407 themselves have a sparse representation in a signal *dictionary* (Zibulevsky et  
 408 al., [62]). Many real-world signals have sparse representations in a proper signal  
 409 dictionary but not in the physical domain. The discussion in Zibulevsky et al. is  
 410 motivated by starting from the case of representing sparse signals in the physical  
 411 domain, depicted in Fig. 8, and then noting that the intuition there carries over  
 412 to the situation of sparsely recovering signals in a transform domain.

## 413 6 Conclusion

414 This paper provided a high-level overview of the philosophy and basic principles  
 415 of the data decomposition approach to data analysis. Starting from the classical  
 416 Singular Value Decomposition method of Linear Algebra and progressing towards  
 417 newer and more powerful methods, such as Independent Component Analysis,  
 418 we showed how the interplay of a geometric depiction of the data space and  
 419 the use of prior constraints on the unknowns can lead to stable solutions to the  
 420 inverse problem of reconstructing the sources. Moreover, we gradually lifted  
 421 the biologically implausible priors imposed by earlier methods and focused on  
 422 the principle of parsimony and on sparsity. These have already given exciting  
 423 results in the field of Computational Neuroscience and promise to give analogous  
 424 results in other fields of Science and Engineering as well.



## 425 Acknowledgements

426 The author is indebted to Prof. Steven Roberts for his invaluable discussions  
427 and stimulating ideas.

## 428 A A Primer on Probability Theory

### 429 A.1 Probability Space

430 The axiomatic formulation of probability starts by defining a probability space,  
431 which is a tuple,  $(\Omega, P)$ , that describes our idea about uncertainty with respect  
432 to a random experiment. It defines:

- 433 • A *sample space*,  $\Omega$ , of possible outcomes,  $\{\omega_i\}$ , of a random experiment  
434 and
- 435 • A probability *measure*,  $P$ , which describes how likely an outcome is.

436 Now, let  $\mathcal{A}$  be a collection of subsets of  $\Omega$ , called random events. Then for  
437  $A \in \mathcal{A}$  the two following conditions must hold:

- 438 • Probabilities must be non-negative,  $P(A) \geq 0$ , and  $P(\Omega) = 1$ ,
- Probabilities must be additive: for two disjoint events,  $A, B$ ,

$$P(A \cap B) = P(A) + P(B) .$$

We also define the *conditional probability*, which can be thought of as “a probability within a probability”,

$$P(A|B) = \frac{P(A \cap B)}{P(B)}, \quad P(B) \neq 0 .$$

Then *random variables* (r.v.’s) are defined as functions from  $\Omega$  to a range,  $\mathcal{R}$ , e.g. a subset of  $\mathbb{R}$  or  $\mathbb{N}$ , etc. These can, inversely, define events as:

$$\mathcal{R} \rightarrow \Omega : \quad A(x) = \left\{ \omega \in \Omega : [x(\omega)] \right\} ,$$

439 where  $[\cdot]$  denotes a “predicate”<sup>8</sup> (e.g. the event ‘ $x > 2$ ’), and therefore act as  
440 “filters” of certain experimental outcomes.

Probability *densities* are defined as densities of probability measures:

$$p(x) = \frac{d}{dx} P(A(x))|_x, \quad \text{with} \quad A(x) = \left\{ x' \in [x, x + dx] \right\}, \quad x \in \mathcal{R} .$$

Finally *joint densities* (e.g. for the case of two random variables  $X, Y$ ) are defined as

$$p_{XY}(x, y) = p \left( \left\{ \omega : X(\omega) = x \wedge Y(\omega) = y \right\} \right) .$$

441 Joint densities of more than two r.v.’s are defined analogously.

<sup>8</sup>This is called an ‘Iverson bracket’ in Iverson notation [35].

## 442 A.2 Three Simple Rules

443 Probability theory is a mathematically elegant theory. The whole construction  
444 can be based on the following three simple rules:

1. The Product rule, which gives the probability of the logical conjunction of two events  $A$  and  $B$ ,

$$P(A \cap B) = P(A|B)P(B) .$$

This can be generalized for  $N$  events, giving the chain rule

$$P\left(\bigcap_{i=1}^N A_i\right) = \prod_{i'=1}^{i-1} P\left(A_i \mid \bigcap_{i'=1}^{i'} A_{i'}\right), \quad i' < i .$$

445 This will be valuable for reasoning in Bayesian networks later.

2. Bayes' rule, which is a recipe that tells us how to update our knowledge in the presence of new information, and can directly be derived from the definition of conditional probability and the product rule,

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}, \quad P(B) \neq 0 .$$

3. Marginalization: given a joint density,  $p_{XY}(x, y)$ , get the marginal density of  $X$  or  $Y$  by integration (i.e. 'integrate out' the uncertainty in one variable):

$$p_X(x) = \int_{\{Y \in \mathcal{Y}\}} p_{XY}(x, y) dy .$$

446 In principle, this is everything we need to know in order to perform proba-  
447 bilistic modelling and inference.

## 448 References

- 449 [1] S.-I. Amari. Natural Gradient Works Efficiently in Learning. *Neural Com-*  
450 *putation*, 10:251–276, 1998.
- 451 [2] H. Asari, B. Pearlmutter, and A. Zador. Sparse representations for the  
452 cocktail party problem. *Journal of Neuroscience*, 26(28):7477–7490, 2006.
- 453 [3] H. Attias. Independent Factor Analysis. *Neural Computation*, 11(4):803–  
454 851, May 1999.
- 455 [4] R. Baraniuk. *Compressive Signal Processing*, 2010.
- 456 [5] H. B. Barlow. Unsupervised Learning. *Neural Computation*, 1(3):295–311,  
457 1989.

- 458 [6] H.B. Barlow, T P Kaushal, and G.J. Mitchison. Finding minimum entropy  
459 codes. *Neural Computation*, 1:412–423, 1989.
- 460 [7] C. Beckmann and S. Smith. Probabilistic independent component analy-  
461 sis for functional magnetic resonance imaging. *IEEE Trans. Med. Imag.*,  
462 23:137–152, 2004.
- 463 [8] A. Bell and T. J. Sejnowski. An information-maximisation approach to  
464 blind separation and blind deconvolution. *Neural Computation*, 7(6):1129–  
465 1159, 1995.
- 466 [9] V.I. Bogachev. *Gaussian measures*. Amer. Math. Soc., 1998.
- 467 [10] V. D. Calhoun, T. Adali, L. K. Hansen, J. Larsen, and J. J. Pekar. ICA of  
468 functional MRI data: An overview, Apr. 2003. Invited Paper.
- 469 [11] J.-F. Cardoso. Infomax and maximum likelihood for blind separation. *IEEE*  
470 *Signal Process. Lett.*, 4(4):112–114, 1997.
- 471 [12] J.-F. Cardoso. Blind signal separation: statistical principles. *Proceedings*  
472 *of the IEEE*, 90(8):2009–2026, Oct. 1998.
- 473 [13] J.-F. Cardoso. The three easy routes to independent component analysis;  
474 contrasts and geometry. In *Proc. of the ICA 2001 workshop*, San Diego,  
475 Dec. 2001.
- 476 [14] Z. Chen and S. Haykin. On different facets of regularization theory. *Neural*  
477 *Computation*, 14:2791–2846, December 2002.
- 478 [15] R. Choudrey, W.D. Penny, and S.J. Roberts. An ensemble learning ap-  
479 proach to independent component analysis. In *Proceedings of the 2000*  
480 *IEEE Signal Processing Society Workshop*, pages 435–444. IEEE Neural  
481 Networks for Signal Processing X, 2000.
- 482 [16] Pierre Comon. Independent component analysis, a new concept? *Signal*  
483 *Processing*, 36(3):287–314, 1994.
- 484 [17] R. T. Cox. Probability, Frequency, and Reasonable Expectation. *Am. Jour.*  
485 *Phys.*, 14:1–13, 1946.
- 486 [18] G. Darmais. Analyse Générale des Liaisons Stochastiques. *Rev. Inst. In-*  
487 *ternat. Stat.*, 21:2–8, 1953.
- 488 [19] I. Daubechies, E. Roussos, S. Takerkart, M. Benharrosh, C. Golden,  
489 K. D’Ardenne, W. Richter, J. Cohen, and J. Haxby. Independent com-  
490 ponent analysis for brain fMRI does not select for independence. *PNAS*,  
491 106(26):10415–10412, 2009.
- 492 [20] A.P. Dempster, N.M. Laird, and D.B Rubin. Maximum Likelihood from  
493 Incomplete Data via the EM Algorithm. *Journal of the Royal Statistical*  
494 *Society. Series B (Methodological)*, 39(1):1–38, 1977.

- 495 [21] D. L. Donoho. Sparse components of images and optimal atomic decom-  
496 positions. *Constructive Approximation*, 17:353–382, 2001.
- 497 [22] D. L. Donoho and M. Elad. Optimally sparse representation in general  
498 (nonorthogonal) dictionaries via  $\ell_1$  minimization. *PNAS*, 100(5):2197–  
499 2202, March 2003.
- 500 [23] D. J. Field. Wavelets, vision and the statistics of natural scenes. *Phil.*  
501 *Trans. R. Soc. Lond. A*, 357:2527–2542, 1999.
- 502 [24] J. H. Friedman and J. W. Tukey. A Projection Pursuit Algorithm for  
503 Exploratory Data Analysis. *IEEE Transactions on Computers*, 23(9):881–  
504 890, Sep. 1974.
- 505 [25] M. Girolami, editor. *Advances in Independent Component Analysis*. Per-  
506 spective on Neural Computing. Springer-Verlag, New York, Aug. 2000.
- 507 [26] M. Girolami. A Variational Method for Learning Sparse and Overcomplete  
508 Representations. *Neural Computation*, 13(11):2517–2532, Nov. 2001.
- 509 [27] C. Golden. *Spatio-Temporal Methods in the Analysis of fMRI Data in*  
510 *Neuroscience*. PhD thesis, Princeton University, 2005.
- 511 [28] Charles F. Golub, Gene H.; Van Loan. *Matrix Computations*. Johns Hop-  
512 kins, 1996.
- 513 [29] T. Hastie, R. Tibshirani, and J. Friedman. *The Elements of Statistical*  
514 *Learning: Data Mining, Inference, and Prediction*. Springer, second edi-  
515 tion, February 2009.
- 516 [30] A. Hyvarinen. The Fixed-Point Algorithm and Maximum Likelihood Es-  
517 timation for Independent Component Analysis. *Neural Processing Letters*,  
518 10(1):1–5, 1999.
- 519 [31] A. Hyvarinen, J. Karhunen, and E. Oja. *Independent Component Analysis*.  
520 Wiley, June 2001.
- 521 [32] A. Hyvarinen and R. Karthikesh. Imposing sparsity on the mixing matrix  
522 in independent component analysis. *Neurocomputing*, 49:151–162, 2002.  
523 Special Issue on ICA and BSS.
- 524 [33] A. Hyvarinen and E. Oja. A Fast Fixed-Point Algorithm for Independent  
525 Component Analysis. *Neural Computation*, 9(7):1483–1492, October 1997.
- 526 [34] A. Hyvarinen and E. Oja. Independent Component Analysis: Algorithms  
527 and Applications. *Neural Networks*, 13(4–5):411–430, 2000.
- 528 [35] Kenneth Iverson. *A Programming Language*. Wiley, 1962.
- 529 [36] D. J. C. Mackay J. W. Miskin. Application of Ensemble Learning ICA to  
530 Infrared Imaging. In *Conference: Independent Component Analysis - ICA*,  
531 2000.

- 532 [37] T. Jaakkola. *Variational Methods for Inference and Estimation in Graphical*  
533 *Models*. PhD thesis, Mass. Inst. of Techn., 1997.
- 534 [38] E.T. Jaynes and G.L. Bretthorst (editor). *Probability Theory: The Logic*  
535 *of Science*. Cambridge University Press, 2003.
- 536 [39] K.H. Knuth. A Bayesian approach to source separation. In C. Jutten  
537 J.-F. Cardoso and P. Loubaton, editors, *Proceedings of the First Interna-*  
538 *tional Workshop on Independent Component Analysis and Signal Separation: ICA'99*, pages 283–288, Aissios, France, Jan 1999.
- 540 [40] H. Lappalainen. Ensemble learning for independent component analysis.  
541 In *Proceedings of the First International Workshop on Independent Com-*  
542 *ponent Analysis*, pages 7–12, 1999.
- 543 [41] H. Lappalainen. Fast Fixed-Point Algorithms for Bayesian Blind Source  
544 Separation, 1999.
- 545 [42] N. D. Lawrence and C. M. Bishop. Variational Bayesian Independent Com-  
546 ponent Analysis. Technical report, , 1999.
- 547 [43] N. Lazar. *The Statistical Analysis of Functional MRI Data*. Statistics for  
548 Biology and Health. Springer, 2008.
- 549 [44] M. S. Lewicki and T. J. Sejnowski. Learning overcomplete representations.  
550 *Neural Computation*, 12(2):337–365, February 2000.
- 551 [45] Y. Li, A. Cichocki, and S.-I. Amari. Analysis of sparse representation and  
552 blind source separation. *Neural Computation*, 16(6):1193–1234, June 1994.
- 553 [46] R. Linsker. A local learning rule that enables information maximization for  
554 arbitrary input distributions. *Neural Computation*, 9:1661–1665, November  
555 1997.
- 556 [47] D. J. C. MacKay. Maximum likelihood and covariant algorithms for inde-  
557 pendent component analysis. Technical report, University of Cambridge,  
558 1996.
- 559 [48] S. G. Mallat and Z. Zhang. Matching Pursuits with Time-Frequency Dic-  
560 tionaries. *IEEE Transactions on Signal Processing*, pages 3397–3415, De-  
561 cember 1993.
- 562 [49] J. W. Miskin and D. J. C. MacKay. Ensemble Learning for Blind Source  
563 Separation. In S. J. Roberts and R. M. Everson, editors, *Independent Com-*  
564 *ponents Analysis: Principles and Practice*. Cambridge University Press,  
565 2001.
- 566 [50] J.-P. Nadal and N. Parga. Nonlinear neurons in the low-noise limit: a  
567 factorial code maximizes information transfer. *Network: Computation in*  
568 *Neural Systems*, 5(4):565–581, 1994.

- 569 [51] B. A. Olshausen and D. J. Field. Emergence of simple-cell receptive field  
570 properties by learning a sparse code for natural images. *Nature*, 381:607–  
571 609, June 1996.
- 572 [52] B. A. Olshausen and D. J. Field. Sparse Coding with an Overcomplete  
573 Basis Set: A Strategy Employed by V1? *Vision Research*, 37(23):3311–  
574 3325, 1997.
- 575 [53] J. A. Palmer, D. P. Wipf, K. Kreutz-delgado, and B. D. Rao. Variational  
576 EM algorithms for non-Gaussian latent variable models. In *Advances in  
577 Neural Information Processing Systems 18*, pages 1059–1066. MIT Press,  
578 2006.
- 579 [54] B. Pearlmutter and L. Parra. Maximum Likelihood Blind Source Sepa-  
580 ration: A context-sensitive generalization of ICA. In *Advances in Neural  
581 Information Processing Systems 9*, pages 613–619, 1997.
- 582 [55] C. R. Rao. A decomposition theorem for vector variables with a linear  
583 structure. *Ann. Math. Statist.*, 40(5):1845–1849, 1969.
- 584 [56] S. Roberts. Independent component analysis: source assessment and sep-  
585 aration, a Bayesian approach. *Vision, Image and Signal Processing, IEE  
586 Proceedings*, 145(3):149–154, Jun. 1998.
- 587 [57] S. Roberts and R. Everson, editors. *Independent Component Analysis:  
588 Principles and Practice*. Cambridge University Press, 2001.
- 589 [58] J. V. Stone. *Independent Component Analysis: A Tutorial Introduction*.  
590 MIT Press, September 2004.
- 591 [59] Y. W. Teh, M. Welling, S. Osindero, G. E. Hinton, T.-W. Lee, J.-F. Car-  
592 doso, E. Oja, and S.-I. Amari. Energy-based models for sparse overcomplete  
593 representations. *Journal of Machine Learning Research*, 4:2003, 2003.
- 594 [60] A. Tonazzini, L. Bedini, and E. Salerno. A Markov Model for Blind Image  
595 Separation by a Mean-Field EM algorithm. *IEEE Trans. Image Proc.*,  
596 15(2), 2006.
- 597 [61] M. Welling and M. Weber. A constrained EM algorithm for Independent  
598 Component Analysis. *Neural Computation*, 13:677–689, 2001.
- 599 [62] M. Zibulevsky and B. A. Pearlmutter. Blind source separation by sparse  
600 decomposition in a signal dictionary. *Neural Computation*, 13:863–882,  
601 April 2001.