

Infomax

The second information theory example we will consider is the Infomax algorithm. This is an approach to unmixing mixed data proposed in Bell and Sejnowski (1995) and gives an in-principle explanation for how the brain might perform auditory source separation. In other words, it shows how source separation might be performed and indicates some of the ingredients that may be present in source separation in the brain.

The problem is as follows: imagine you are in a crowded room, in the classic telling, at a cocktail party. Lots of people are talking, but if you concentrate on one voice you can separate it from the overall hub bub. We can experience the opposite effect when meditating or sitting in thoughtful silence: we suddenly notice how loud distance noises are, the noise of voices on the street, other people coughing and clearing their throats, the sound of wind through nearby trees. All these are sounds we automatically filter out when listening. This process of picking out individual sources of sound from a mixture of sounds is called auditory source separation and the question of how to do this is called the cocktail party problem. In fact, the problem is much more general than auditory source separation; the same techniques should apply to any problem where we can observe a mixed signal and would like the unmixed signal. In this more general context the problem of source separation is called independent component analysis.

As noted above, we are interested in the particularly neuronal approach described in Bell and Sejnowski (1995); however this isn't the only approach to the problem and a different algorithm, fastICA (Hyvarinen, 1999), is the most popular algorithm for source separation.

The problem can be phrased like this, given the sources $\mathbf{s}(t)$ where \mathbf{s} is a vector over multiple sources. Now we do source separation with only two recordings, one for each ear; here we are just going to consider the simpler problem of source separation when there are as many recordings as there are sources, we also assume the mixing is linear and instantaneous, real mixing of auditory signals in a room will only have these properties approximately. However, given these assumption we have

$$\mathbf{r}(t) = M\mathbf{s}(t) \tag{1}$$

with M a square mixing matrix. The goal is to find the unknown source signals $\mathbf{s}(t)$ from the known recordings $\mathbf{r}(t)$. Although it makes little difference,

let's restrict ourselves to two sources, so \mathbf{s} and \mathbf{r} are two-dimensional vectors and M is a two-by-two matrix.

We assume the two sources are independent, $p_{S_1, S_2} = p_{S_1} p_{S_2}$; the point of source separation is to unmix independent sources. Now, we want to find an unmixing matrix W so that knowing

$$\mathbf{x}(t) = W\mathbf{r}(t) \quad (2)$$

is as good as knowing the sources; precisely, since $\mathbf{x} = WM\mathbf{s}$ we want WM to be a diagonal matrix multiplied by a permutation matrix, we do not mind if unmixing changes the overall amplitude of the source, or if it reorders the sources. In this two-to-two example, that means

$$MW = \text{diag}(d_1, d_2) \quad (3)$$

or

$$MW = \begin{pmatrix} 0 & d_1 \\ d_2 & 0 \end{pmatrix} \quad (4)$$

where d_1 and d_2 are real numbers. Hence:

$$\mathbf{s} \xrightarrow{\text{mixing}} \mathbf{r} = M\mathbf{s} \xrightarrow{\text{unmixing}} \mathbf{x} = W\mathbf{r} \quad (5)$$

One difficulty with looking at this problem is that it involves continuous random variables, whereas the unit so far has concentrated on the discrete case: $s_1(t)$ and $s_2(t)$ are continuous variable. For this reason some of the results will just be quoted and any exam questions will reflect the difficulty of this topic. Recall that the information is defined in the same way as before

$$h(X) = - \int p(x) \log p(x) dx \quad (6)$$

but it is no longer always positive. Furthermore, we saw that

$$h(\lambda X) = h(X) + \log |\lambda| \quad (7)$$

Now, the idea is to solve the problem by using the fact that S_1 and S_2 are independent: we just need to find W so that X_1 and X_2 are also independent. One approach might be to decorrelate the random variables:

$$C(X_1, X_2) = \langle (X_1 - \langle X_1 \rangle)(X_2 - \langle X_2 \rangle) \rangle_{(X_1, X_2)} \quad (8)$$

where the expectation value for continuous random variables has the obvious definition

$$\langle g(X) \rangle = \int p_X(x) g(x) dx \quad (9)$$

It is easy to check that the correlation vanishes if X_1 and X_2 are independent, however, the flaw in this approach is that the converse is not true, it is possible to have zero correlation while still having statistical dependence. To see this, imagine, for convenience and without loss of generality, that EX_1 and EX_2 are zero and that we have chosen W so that the correlation matrix is the identity:

$$C_{ab} = C(X_a, X_b) = \mathbf{1} \quad (10)$$

then it is easy to see that rotations

$$\begin{pmatrix} X'_1 \\ X'_2 \end{pmatrix} = \begin{pmatrix} \cos \theta & \sin \theta \\ -\sin \theta & \cos \theta \end{pmatrix} \begin{pmatrix} X_1 \\ X_2 \end{pmatrix} \quad (11)$$

do not change the correlation matrix. Thus, the decorrelation prescription has a rotational ambiguity and something more is needed. That something, of course, is to require $I(X_1, X_2) = 0$ since this happens if and only if X_1 and X_2 are independent.

The problem is that $I(X_1, X_2)$ is difficult to calculate: the idea behind infomax is to look at $h(X_1, X_2)$:

$$I(X_1, X_2) = h(X_1) + h(X_2) - h(X_1, X_2) \quad (12)$$

The idea is that maximizing the joint entropy, $h(X_1, X_2)$ will give a minimum of the mutual information, in other words, the variations in the individual entropies $h(X_1)$ and $h(X_2)$ can be ignored. However, as stated, this will not work because the entropy can be increased by a trivial scaling, $X_a \rightarrow \lambda X_a$, changes the joint entropy, $h(X_1, X_2) \rightarrow h(X_1, X_2) + \log |\lambda|$ so $h(X_1, X_2)$ can be made arbitrarily large by scaling, something that tells us nothing about mixing and unmixing. Inspired by the behaviour of neurons, in Bell and Sejnowski (1995) this is solved by adding a saturation non-linearity:

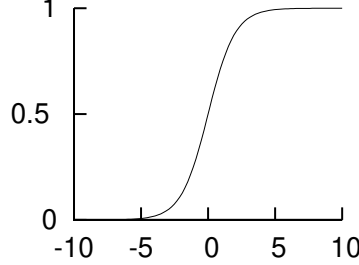
$$\begin{aligned} y_1 &= g(x_1 + w_1) \\ y_2 &= g(x_2 + w_2) \end{aligned} \quad (13)$$

where w_1 and w_2 are parameters and, for example,

$$g(u) = \frac{1}{1 + e^{-u}} \quad (14)$$

is a saturating non-linearity so $g : (-\infty, \infty) \rightarrow (0, 1)$.

(15)



Now we have¹

$$\mathbf{s} \xrightarrow{\text{mixing}} \mathbf{r} = M\mathbf{s} \xrightarrow{\text{unmixing}} \mathbf{x} = W\mathbf{r} \xrightarrow{\text{non-linearity}} \mathbf{y} : y_a = g(x_a + w_a) \quad (16)$$

For later notational convenience, let's write

$$y_a = g(x_a + w_a) = f(r_1, r_2; W, w_a) \quad (17)$$

where f is the function, parameterized by W and w_a , mapping from the recording to y . The reason the saturating non-linearity will help is that with a non-linearity like this a large enough scaling will reduce the entropy: consider $h(g(\lambda X))$ where λ is a constant. If λ is very large it will spread X out so that it will take lots of very big positive or negative values; if λx is a big positive number then $g(\lambda x)$ will be near to one, conversely, if it is a large negative number, $g(\lambda x)$ will be near to zero: a distribution that is often zero and often one is not very entropic.

Before considering the source separation problem, let's look at the effect of the non-linearity on its own: we consider the one-to-one case

$$r \xrightarrow{\text{multiply}} x = Wr \xrightarrow{\text{non-linearity}} y = g(x + w) = f(r; w, W) \quad (18)$$

where W and w are now both scalars and r , x and y are outcomes for random variables R , X and Y . We consider maximizing the entropy $h(Y)$. What does this do; well, it maximizes the information in Y about R :

$$I(R; Y) = h(Y) - h(Y|R) \quad (19)$$

¹In broadcast notation $\mathbf{y} = g.(\mathbf{x} + \mathbf{w})$.

but $h(Y|R)$ is constant since R determines Y . In the discrete case we are familiar with this would be easy to discuss, $h(Y|R)$ would be zero, in the continuous case it is not that simple, it is actually minus infinity, but, the consequence is the same, it does not depend on W and w . To maximize $h(Y)$ we need to calculate its derivative with respect to the parameters W and w ; this is a feature of the algorithm, ultimately we need to calculate derivatives and these are calculable and well defined, even if the quantity being differentiated is not. In this case

$$h(Y) = - \int p(y) \log p(y) dy \quad (20)$$

and this is estimated by

$$\tilde{h}(y) = -\log p(y) \quad (21)$$

In other words, if n values of y_i are drawn from Y then

$$\frac{1}{n} \sum_i \tilde{h}(y_i) \rightarrow h(Y) \quad (22)$$

as n gets large.

Of course we do not have $p_Y(y)$ and it would be difficult to estimate, but, it turns out we do not need it to get the derivative of $\tilde{h}(y)$. We have seen already that since $y = f(r; W, w)$

$$p_Y(y) = \frac{p_R[r = f^{-1}(y)]}{|f'(f^{-1}(y))|} \quad (23)$$

so

$$\tilde{h}(y) = -\log p_R(r) + \log |f'| \quad (24)$$

and $p_R(r)$ is independent of the parameters. Now, for our choice of saturating non-linearity

$$\begin{aligned} g(u) &= \frac{1}{1 + \exp(-u)} \\ \frac{dg}{du} &= g(1 - g) \end{aligned} \quad (25)$$

and hence

$$\log |f'| = \log W + \log f + \log (1 - f) \quad (26)$$

Now we know f :

$$f = g(Wr + w) \quad (27)$$

so

$$\frac{df}{dW} = rf(1 - f) \quad (28)$$

and hence,

$$\frac{d\tilde{h}(y)}{dW} = \frac{1}{W} + \frac{1}{f}rf(1 - f) - \frac{1}{1 - f}rf(1 - f) = \frac{1}{W} + r(1 - 2y) \quad (29)$$

Similarly

$$\frac{d\tilde{h}(y)}{dw} = 1 - 2y \quad (30)$$

These quantities: s , y and, of course, W , are numbers we have access to, W is a parameter, s is the signal, we can sample $s(t)$ at a set of times to get a set of s 's and y is a function of s . This means we can estimate these derivatives, giving the gradient of h at a point (W, w) in the parameter space. This is a common situation in numerical optimization, we don't know the function and, here, it isn't even so easy to define, but we do know its gradient. This means that locally we know which direction brings us in the direction of greater h and so numerical hill-climbing routines can be used, these are well described in, for example, Press et al. (2007): steepest ascent, conjugate gradient or metric gradient methods work well here.²

The idea then is to choose a starting W and w ; estimate the gradient and then change W and w a small amount, repeating until the optimum values are found. What would the optimum value look like, well we know

$$p_Y(y) = \frac{p_R(r)}{|f'(r)|} \quad (31)$$

with $r = f^{-1}(y)$. Hence, if Y is evenly distributed on its interval $(0, 1)$ and $f'(r)$ is always positive then

$$f(r) = \int_{-\infty}^r p_R(u) du \quad (32)$$

²The metric method uses a slightly surprising, but effective, choice of metric and is defined in Amari (1998).

Now this is not what we do, we do not know $p_R(r)$ and we chose $f(r)$ at the start, here it's a member of a two-parameter family of functions parameterized by W and w . However, ideally, if the derivative of the saturating non-linearity is somewhat close to the distribution of R then infomax will find the W and w that line everything up so that Y will have something close to an even distribution. Here is an example:

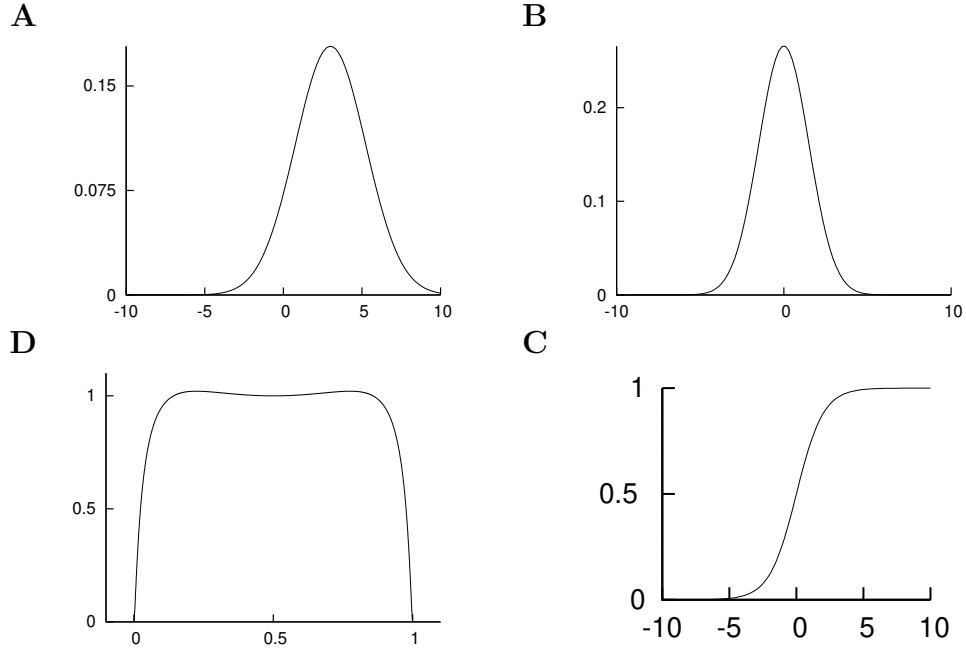


Figure **A** shows an initial distribution

$$p_R(r) = \frac{1}{\sqrt{10\pi}} e^{-(r-3)^2/10} \quad (33)$$

B is a new distribution arrived at by shifting and rescaling r using W and w :

$$p_U(u) = \frac{1}{\sqrt{4.5\pi}} e^{-u^2/4.5} \quad (34)$$

where $u = Wr + w$. Now, this distribution is all lined up with the rectifying non-linearity **C** so that $p_Y(y)$ where $y = g(u)$ is **D**.

Now, back to the two-to-two case:

$$\mathbf{s} \xrightarrow{\text{mixing}} \mathbf{r} = M\mathbf{s} \xrightarrow{\text{unmixing}} \mathbf{x} = W\mathbf{r} \xrightarrow{\text{non-linearity}} \mathbf{y} : y_a = g(x_a + w_a) \quad (35)$$

and here we want to maximize $h(Y_1, Y_2)$; the idea being that this should find a matrix W whose eigen-directions give statistically independent Y_a , this is the bit we want since it will also make the X_a independent, and whose eigenvalues, along with the values of w_a make $h(Y_1)$ big by lining the saturating non-linearity up with the underlying distributions: the orientation of W deals with the unmixing, the scale of W and the vector of w_a s deals with the one-to-one part. Anyway, doing the calculation gives

$$\begin{aligned}\frac{d\tilde{h}(y)}{dW_{ab}} &= (W^T)_{ab}^{-1} + r_a(1 - 2y_b) \\ \frac{d\tilde{h}(y)}{dw_a} &= 1 - 2y_a\end{aligned}\tag{36}$$

allowing the maximum of $h(Y_1, Y_2)$ to be, hopefully, found and this, again, hopefully, will unmix the signal. Note that this algorithm is not as blind as we might of hoped, the non-linearity needs to be chosen judiciously: however, the algorithm is reasonably robust; reasonably successful and certainly interesting mathematically.

References

- Amari, S.-I. (1998). Natural gradient works efficiently in learning. *Neural computation*, 10(2):251–276.
- Bell, A. J. and Sejnowski, T. J. (1995). An information-maximization approach to blind separation and blind deconvolution. *Neural computation*, 7(6):1129–1159.
- Hyvarinen, A. (1999). Fast and robust fixed-point algorithms for independent component analysis. *IEEE transactions on Neural Networks*, 10(3):626–634.
- Press, W. H., Teukolsky, S. A., Vetterling, W. T., and Flannery, B. P. (2007). *Numerical recipes 3rd edition: The art of scientific computing*. Cambridge university press.