# Blind Source Separation 2.

We first need to introduce the idea of data "whitening". First, given a vector of data $\vec{x}(K)$, we can always take

$$\vec{x}(K) - E(\vec{x}_k)$$

to transform it into zero-mean data. We'll apply this without further mention to assume that <u>all</u> our data is ~~previously~~ zero mean.

Given such a data vector $\vec{x}(K)$, we let

$$R_{xx} = E(xx^T)$$

be the covariance matrix, as usual.

If we take the SVD of $R_{xx}$, we get

$$R_{xx} = V_x \Lambda_x V_x^T$$

We can now apply $V_x$ to the data to get a new time series:

$$\vec{y}(k) = V_x^T \vec{x}(k)$$

Notice that

$$\vec{y}\,\vec{y}^T = V_x^T \vec{x}\,\vec{x}^T V_x$$

So

$$E(\vec{y}\,\vec{y}^T) = V_x^T E(\vec{x}\,\vec{x}^T) V_x$$
$$= V_x^T (V_x \Lambda_x V_x^T) V_x$$
$$= \Lambda_x,$$

so the entries in $\vec{y}(k)$ are uncorrelated.

This process, reasonably enough, is called <u>decorrelation</u>. On the other hand, the entries in $\vec{y}(k)$ <u>don't</u> have unit variance. We solve that by rescaling.

Lemma. If $\vec{x}(k)$ is zero-mean and

$$R_{xx} = E(\vec{x}\vec{x}^T) = V_x \Lambda_x V_x^T, \text{ then}$$

$$\vec{y}(k) = \Lambda_x^{-1/2} V_x^T \vec{x}(k)$$

has $R_{yy} = E(\vec{y}\vec{y}^T) = I_n$.

Proof. We just write out

$$R_{yy} = E(\Lambda_x^{-1/2} V_x^T \vec{x}\vec{x}^T V_x \Lambda_x^{-1/2})$$

$$= \Lambda_x^{-1/2} V_x^T ~\rlap{\text{\scriptsize \sout{XXXX}}}\phantom{XXX}~ R_{xx} V_x \Lambda_x^{-1/2}$$

$$= \Lambda_x^{-1/2} V_x^T (V_x \Lambda_x V_x^T) V_x \Lambda_x^{-1/2}$$

$$= \Lambda_x^{-1/2} \Lambda_x \Lambda_x^{-1/2} = I_n. \qquad \square$$

We call this a "whitening" transformation because

Definition. A random vector $\check{X}(K)$ is called a <u>white noise vector</u> if the components of $\check{X}(K)$ are selected from <u>statistically independent</u> <u>probability distributions</u> <u>with</u> <u>zero</u> <u>mean</u> <u>and</u> <u>finite variance</u>.

Now suppose we have our usual BSS setup

$$\vec{X}(K) = H \vec{s}(K)$$

If we ~~tran~~ whiten the data, we get

$$\vec{y}(K) = Q \vec{X}(K) = QH \vec{s}(K).$$

(where $Q = \Lambda_x^{-1/2} V_x^T$, as above.)

Notice that if we take $A = QH$ as the new mixing matrix, we have

$$\vec{y}(K) = A\,\vec{s}(K),$$

so

$$R_{yy} = E(\vec{y}\vec{y}^T)$$

$$= E(A\vec{s}\vec{s}^T A^T)$$

$$= A\,E(\vec{s}\vec{s}^T)\,A^T.$$

But we assumed that $E(\vec{s}\vec{s}^T) = R_{ss} = I_n$ and we just assured that $R_{yy} = I_n$, so this implies

$$I_n = AA^T$$

and the new mixing matrix is <u>orthogonal</u>. Thus the <u>demixing</u> matrix $W$ should be given by $A^T = A^{-1} = W$.

Now suppose we only want to get e of the n signals back. We'd like to express the problem in terms of minimizing the _mutual information_ of the ~~to~~ recovered signals.

The mutual information is given by Kullback-Leibler divergence.

The Kullback-Leibler divergence is a measure of the difference between probability measures.

If we have measures $P, Q$ ~~$M$~~ on $X$ so that

~~$M$~~ $P$ is absolutely continuous with respect to $Q$

$-$ or $-$

~~$M$~~ $Q(A) = 0 \Rightarrow P(A) = 0.$

~~$M$~~ ~~$The$~~

Then the Radon-Nikodym derivative $\frac{dP}{dQ}$ exists. ~~and is~~ This is a measurable function so that

$$P(A) = \int_{a \in A} \left(\frac{dP}{dQ}\right)(a) \, dQ$$

for all subsets ~~of~~ $A$ of $X$.

we can then write

$$D_{KL}(P \| Q) = \int_X \left( \frac{dP}{dQ} \ln\left(\frac{dP}{dQ}\right) \right) dQ$$

which is the entropy of P with respect to Q. If Q and P are the same measure, then this is zero, otherwise it increases as P differs from Q by larger and larger amounts.

We can use this idea to measure whether two ~~distributions~~ variables are independent. Here's the idea: ~~If P and Q are independent, then on X×X~~

JH

Suppose we are given a measure $\mu$ on $X \times X$. We can generate two ~~distributions~~ measures on $X$ by pushing forward by

$$\pi_1 : X \times X \to X, \quad \pi_1(x,y) = x$$

$$\pi_2 : X \times X \to X, \quad \pi_2(x,y) = y$$

to two new measures $(\pi_1)_* \mu$, $(\pi_2)_* \mu$.

If $\mu$ is the joint distribution of independent measures on $X$, then

$$\mu = (\pi_1)_* \mu \times (\pi_2)_* \mu,$$

$$\underset{\text{\underline{marginal} distributions}}{\nwarrow \quad \uparrow}$$

so it makes sense to try to measure the "degree of independence" of $x$ and $y$ by computing KL-divergence from the

joint distribution to the product of the marginals.

Definition. The <u>mutual information</u> of random variables $(X_1, X_2)$ sampled according to $\nu$ on $X \times X$ is defined by

$$I(X_1, X_2) = D_{KL}\left(\nu \,\|\, (\pi_1)_* \nu \cdot (\pi_2)_* \nu\right).$$

~~We can view BSS as the problem of choosing a demixing matrix We so that the product (summed)~~

We can extend this definition to measure the mutual information of any number of variables by

$$I(X_1, \ldots, X_K) = D_{KL}\left(\nu \,\|\, (\pi_1)_* \nu \times \cdots \times (\pi_K)_* \nu\right).$$

So suppose we have (prewhitened) signals $\vec{X}(\ell)$. Amari's minimization of MI algorithm proposes that we solve

$$\text{"} \quad \underset{W_K \in V_K(\mathbb{R}^n)}{\text{minimize}} \; I\left(W_K \vec{X}(\ell)\right) \text{"}.$$

Of course, there is a serious question here: how do we estimate the entropy of $\nu$ with respect to $(\Pi_1)_* \nu \times \dots \times (\Pi_K)_* \nu$ from the data?

~~We first note~~

$$D_{KL}(P\|Q) = -E_p\left(\ln q(x)\right) + E_p\left(\ln p(x)\right)$$

Amari observes that <u>when the number of measurements n is equal to the number of unmixed signals K</u> ~~the~~

~~entropizing~~
~~∄ H The KL divergence can be written as~~

~~𝔼~~

if we write the joint distribution of the signals as $\nu$ (a function on $\mathbb{R}^n$) ~~and the~~ wrt Lesbègue measure on $\mathbb{R}^n$ and the marginal distributions as $\nu_i$ ~~#~~ (the density of $(\pi_i)_* \nu$ on $\mathbb{R}$ ~~as~~ with respect to Lesbègue measure),

then

$$I(\vec{X}) = \int_{\mathbb{R}^n} \nu \log \frac{\nu}{\Pi \nu_i} \, dVol$$

$$= \int_{\mathbb{R}^n} \nu \log \nu \, dVol - \int_{\mathbb{R}^n} \nu \log \Pi \nu_i \, dVol$$

$$= \int_{\mathbb{R}^n} \nu \log \nu \, dVol - \int_{\mathbb{R}^n} \nu \log \nu_i \, dVol.$$

The first term is the "differential entropy" of $\nu$. Using the change of variables formula, if we transform by another mixing matrix $W$, we get the differential entropy to.

change by writing the new density
as $(\det W)\nu$. In this case, we get

$$\int (\det W \cdot \nu) \log (\det W \cdot \nu) \overset{(\det W)^{-1}}{dVol}$$

$$= \int \nu \log (\det W \cdot \nu) dVol$$

$$= \int \nu \log \nu \, dVol + \int \nu \log (\det W) \, dVol$$

$$= \int \nu \log \nu \, dVol + \log (\det W).$$

or "old entropy + $\log (\det W)$ = new entropy"

Now we have to consider

$$\int_{\mathbb{R}^n} \nu \log \nu_i \, dVol.$$

Now

$$\nu_i(x_i) = \int_{\mathbb{R}^{n-1}} \nu(x_1,\ldots,x,\ldots x_n)\, dx_1 \wedge \ldots \wedge \widehat{dx_i} \wedge \ldots \wedge dx_n$$

so we have $\nu_i$ only depending on $x_i$.

Thus

$$\int_{\mathbb{R}^n} \nu \log \nu_i \, dx_1 \wedge \ldots \wedge dx_n =$$

$$\int_{x_i} \log \nu_i \left( \int_{\mathbb{R}^{n-1}} \nu(x_1,\ldots,x,\ldots,x_n)\, dx_1 \wedge \ldots \wedge \widehat{dx_i} \wedge dx_n \right) dx_i$$

$$= \int_{x_i} \nu_i \log \nu_i \, dx_i,$$

which is just the (differential) entropy of the $i$th (demixed) signal.

We now face a challenge:

1) The joint entropy doesn't depend on demixing matrix $W_n$.

2) The entropy of product of marginals with respect to joint density is simply a sum of one-variable entropies.

Now we must estimate _these_ entropies, from the data.