

Maximizing mutual information between input and neural response

S. Marzen

March 9, 2021

This is basically ICA in a certain limit with added noise.

We will deal with the one-dimensional case first:

$$z = f(Ax) \tag{1}$$

$$n = z + B\eta \tag{2}$$

where η is Gaussian noise, but we assume that such noise is small enough to ignore, and where z is the mean firing rate of the neuron. Then we have

$$I[x; n] = H[n] - H[n|x] \tag{3}$$

$$= H[n] + \frac{1}{2} \log(2\pi e B^2) \tag{4}$$

and with

$$\rho(n) = \int_{-\infty}^{\infty} \rho_z(n-y) \times \frac{1}{\sqrt{2\pi B^2}} e^{-y^2/2B^2} dy. \tag{5}$$

Since B is small, we simply Taylor expand and find

$$\rho_z(n-y) = \rho_z(n) - \frac{\partial \rho_z}{\partial z} \Big|_n y + \frac{1}{2} \frac{\partial^2 \rho_z}{\partial z^2} \Big|_n y^2 + O(y^3) \tag{6}$$

and thus

$$\rho(n) = \rho_z(n) + \frac{1}{2} \rho_z''(n) B^2 + O(B^4). \tag{7}$$

To find its differential entropy, we simply have

$$H[n] = - \int \rho(n) \log \rho(n) dn \tag{8}$$

$$= - \int \left(\rho_z(n) + \frac{1}{2} \rho_z''(n) B^2 \right) \log \left(\rho_z(n) + \frac{1}{2} \rho_z''(n) B^2 \right) + O(B^4) \tag{9}$$

$$= - \int \rho_z(n) \log \rho_z(n) dn - \frac{1}{2} B^2 \int \rho_z''(n) \log \rho_z(n) dn - \frac{1}{2} B^2 \int \rho_z''(n) dn + O(B^4). \tag{10}$$

Although we have the noise terms in there, we will simply ignore them in order to understand the predictions that come from maximizing mutual information between stimulus and neural response. We have (and this is in Cover and Thomas)

$$\rho_z(z) = \frac{\rho_x(x)}{|dz/dx|} \tag{11}$$

$$H[n] = - \int \rho_z(z) \log \frac{\rho_x(x)}{|dz/dx|} dz + O(B^2) \tag{12}$$

$$= - \int \rho_x(x) \log \rho_x(x) + \int \rho_x(x) \left| \frac{dz}{dx} \right| dx + O(B^2) \tag{13}$$

$$= H[x] + \int \rho_x(x) \log \left| \frac{dz}{dx} \right| dx + O(B^2). \tag{14}$$

We make yet another small noise approximation— that x is highly peaked at some value x^* and can be approximated as a Gaussian around that point.

$$\frac{dz}{dx} = Af'(Ax) \quad (15)$$

$$= Af'(Ax^*) + A^2 f''(Ax^*)(x - x^*) + \frac{1}{2} A^3 f'''(Ax^*)(x - x^*)^2 + O((x - x^*)^3) \quad (16)$$

$$\log \left| \frac{dz}{dx} \right| = \log |Af'(Ax^*) + A^2 f''(Ax^*)(x - x^*) + \frac{1}{2} A^3 f'''(Ax^*)(x - x^*)^2 + O((x - x^*)^3)| \quad (17)$$

$$= \log |Af'(Ax^*)| + \frac{Af''(Ax^*)}{f'(Ax^*)}(x - x^*) + \frac{A^2 f'''(Ax^*)}{2f'(Ax^*)}(x - x^*)^2 + O((x - x^*)^3) \quad (18)$$

and so

$$H[n] = H[x] + \int \rho_x(x) \left(\log |Af'(Ax^*)| + \frac{Af''(Ax^*)}{f'(Ax^*)}(x - x^*) + \frac{A^2 f'''(Ax^*)}{2f'(Ax^*)}(x - x^*)^2 + O((x - x^*)^3) \right) dx + O(B^2) \quad (19)$$

$$= H[x] + \log |A| + \log |f'(Ax^*)| + \frac{A^2 f'''(Ax^*)}{2 f'(Ax^*)} \sigma^2 + O(\sigma^4) \quad (20)$$

where σ is $\langle (x - x^*)^2 \rangle$. Thus,

$$I[x; n] = H[x] + \frac{1}{2} \log(2\pi e B^2) + \log |A| + \log |f'(Ax^*)| + \frac{A^2 f'''(Ax^*)}{2 f'(Ax^*)} \sigma^2 + O(B^2, \sigma^4). \quad (21)$$

To maximize this, we want to make A as large as possible in magnitude, we want $|f'(Ax^*)|$ to be as close to zero as possible but positive, and we want $f'''(Ax^*)$ to be as large as possible. The last of these leads to the weird conclusion that if we have a maximum in firing rate at x^* , then $f'''(Ax^*)$ being large and positive will cause the downturn in firing rate to swing into an upturn. The firing rate might only have a local max at the most likely value of the input.

One can do this for the multi-dimensional case. For ease (although this can be done more generally) we assume that the dimensionality of input and neural responses match. Then, given that

$$n = \vec{f}(Ax) + B\eta \quad (22)$$

where A is now a matrix and η is now a vector of ignorable noise, one finds that

$$I[x; n] = \log \left| \det \left(A^\top \text{diag}(\vec{f}'(x^*)) \right) \right| + \frac{1}{2} \text{tr} \left(\text{diag}(\vec{f}'(x^*)^{\odot -1}) A^{-\top} \text{diag}(\vec{f}'''(x^*)) A^\top \text{diag} \left((A \odot A) \sigma^{\vec{2}} \right) \right),$$

where x^* is the most likely input and $\sigma^{\vec{2}}$ is $\langle (x - x^*)(x - x^*)^\top \rangle$. The qualitative conclusions from the single neuron case seem to hold up.

The derivation goes as follows. We again ignore to $O(B^2)$ and to $O(\sigma^4)$ and find that, as is typical for ICA, the only thing that matters is $\langle \log |\det \nabla_x \vec{f}| \rangle$. The strategy again is to Taylor expand $\nabla_x \vec{f}$ around x^* :

$$\left(\nabla_x \vec{f} \right)_{i,j} = \frac{\partial f_j}{\partial x_i} \quad (23)$$

$$= \frac{\partial f_j}{\partial x_i} \Big|_{x^*} + \sum_k \frac{\partial^2 f_j}{\partial x_i \partial x_k} \Big|_{x^*} (x_k - x_k^*) + \frac{1}{2} \sum_{k,k'} \frac{\partial^3 f_j}{\partial x_i \partial x_k \partial x_{k'}} \Big|_{x^*} (x_k - x_k^*)(x_{k'} - x_{k'}^*) \quad (24)$$

and so we can write

$$\nabla_x \vec{f} = M(x^*) + \delta M \quad (25)$$

where δM is “small” in a way to be described. Then

$$\langle \log |\det(M(x^*) + \delta M)| \rangle = \langle \log |\det M(x^*)| \rangle + \langle \text{tr}(M(x^*)^{-1} \delta M) \rangle. \quad (26)$$

We have

$$\langle \log |\det M(x^*)| \rangle = \log |\det M(x^*)| \quad (27)$$

$$= \log |\det(A^\top \text{diag}(\vec{f}'(Ax^*)))| \quad (28)$$

$$= \log |\det A| + \sum_i \log f'_i(Ax^*) \quad (29)$$

and then we have

$$\langle \text{tr}(M(x^*)^{-1} \delta M) \rangle = \text{tr}(M(x^*)^{-1} \langle \delta M \rangle) \quad (30)$$

with

$$\langle \delta M_{ij} \rangle = \frac{1}{2} \sum_{k,k'} \frac{\partial^3 f_j}{\partial x_i \partial x_k \partial x_{k'}} \Big|_{x^*} \delta_{k,k'} \sigma_k^2 \quad (31)$$

$$= \frac{1}{2} \sum_k \frac{\partial^3 f_j}{\partial x_i \partial x_k \partial x_k} \Big|_{x^*} \sigma_k^2 \quad (32)$$

where the components of x are assumed to be independent. (It is rather easy to transform your data so that this is roughly true around the peak.) This is unfortunately ugly, but hopefully useful.