

# APPENDIX

## A ON THE UNIDENTIFIABILITY OF NONLINEAR ICA

The purpose of this section is to briefly review the proof of unidentifiability of nonlinear ICA as [22]: In this section we assume the most general conventional form of nonlinear ICA where the generative model follows:

$$\mathbf{x} = \mathbf{f}(s) \quad (22)$$

where  $s$  are the independent sources and  $\mathbf{x}$  are mixed signals. In the following, we show how to construct a function  $\mathbf{g} : R^n \rightarrow R^n$  so that the components  $\mathbf{y} = \mathbf{g}(\mathbf{x})$  are independent. More importantly, we show that this construction is by no means unique.

### A.1 EXISTENCE

The proposed method in [22] is a generalization of the famous Gram-Schmidt orthogonalization. Given  $m$  independent variables,  $y_1, \dots, y_m$  and a variable  $x$ , one constructs a new variable  $y_{m+1} = g(y_1, \dots, y_m, x)$  so that the set  $y_1, \dots, y_{m+1}$  is mutually independent. The construction process is defined recursively as follows. Assume we have  $m$  independent random variables  $y_1, \dots, y_m$  with uniform distribution in  $[0, 1]^m$ .  $x$  is any random variable and  $a_1, \dots, a_m, b$  are some nonrandom scalars. Next, we define

$$\begin{aligned} g(a_1, \dots, a_m, b; p_{y,x}) &= p(x \leq b | y_1 = a_1, \dots, y_m = a_m) \\ &= \frac{\int_{-\infty}^b p_{y,x}(a_1, \dots, a_m, \xi) d\xi}{p_y(a_1, \dots, a_m)} \end{aligned} \quad (23)$$

Theorem 1 of [22] says that the random variable defined as  $y_{m+1} = g(y_1, \dots, y_m, x)$  is independent from the  $y_1, \dots, y_m$  and  $y_1, \dots, y_{m+1}$  are uniformly distributed in the unit cube  $[0, 1]^{m+1}$ .

### A.2 NON-UNIQUENESS

In the previous section, it was shown that there exists a mapping  $\mathbf{g}$  that transforms any random vector  $\mathbf{x}$  into a uniformly distributed random vector  $\mathbf{y} = \mathbf{g}(\mathbf{x})$ . Here, we show that the construction of  $\mathbf{g}$  is not unique and this non-Uniqueness can be caused by several factors.

- A linear transformation  $\mathbf{x}'$  can precede the nonlinear map  $\mathbf{f}$  and then compute the independent components  $\mathbf{y}' = \mathbf{g}'(\mathbf{x}')$  where  $\mathbf{g}'$  is computed as describe in the previous section. The new map  $\mathbf{g}'$  gives

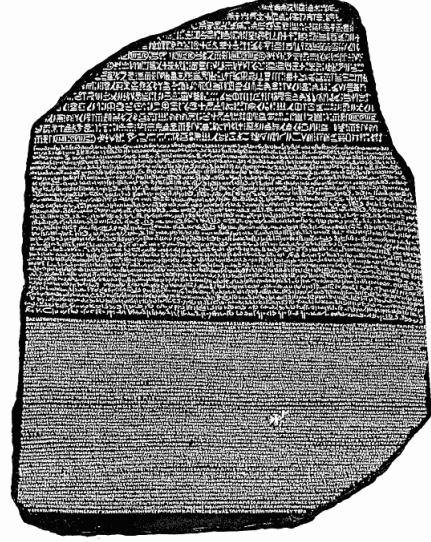


Figure 4: The Rosetta Stone, a stele found in 1799, inscribed with three versions of a decree issued at Memphis, Egypt in 196 BC. The top and middle texts are in Ancient Egyptian using hieroglyphic script and Demotic script, respectively, while the bottom is in Ancient Greek. (Source: Wikipedia)

a new decomposition of  $\mathbf{x}$  into independent components  $\mathbf{y}'$  which can not be trivially reduced to  $\mathbf{y}$ .

- An element-wise function  $\mathbf{h}$  can apply on the independent sources  $s$  first to give new sources  $s'$  such that  $s'_i = h_i(s_i)$ . Constructing the solution  $g$  for these new scaled version of sources gives a new decomposition into independent components.
- Assume a class of measure-preserving automorphisms  $\mathbf{h} : [0, 1]^n \rightarrow [0, 1]^n$ . The mapping  $\mathbf{h}$  does not change the probability distribution of a uniformly distributed random variable in  $n$ -dimensional hypercube. The composition  $\mathbf{h} \circ \mathbf{g}$  gives another solution to nonlinear ICA. Therefore, the class of measure-preserving automorphisms gives a parameterization of the solutions to nonlinear ICA introducing a class of non-trivial indeterminacies.

If only independence among the components matters, it is possible to construct a mapping  $\mathbf{y} = G(\mathbf{x})$  such that  $y_i$  is independent of  $y_j$  for  $i \neq j$  and uniformly distributed in  $[0, 1]^n$ . This shows that at least one solution exists. The non-uniqueness of the solution can be shown by parameterising a class of infinitely many solutions. Once  $\mathbf{y}$  is found with above conditions, any measure-preserving automorphism  $f : [0, 1]^n \rightarrow [0, 1]^n$  can be used to parameterize  $G$  as  $f \circ G$ , suggesting that there

are infinitely many solutions to nonlinear ICA whose relations are nontrivial.

### A.3 THE SCALAR INVERTIBLE FUNCTION GAUGE

Another indeterminacy is element-wise functions  $f_i$  applying on  $y_i$  which suggests another dimension of ambiguity. Non-Gaussianity cannot help here since we can construct any marginal distribution by combining the CDF of the observed variable with the inverse CDF of the target marginal distribution. This indeterminacy is in some sense unavoidable and is related to the fact that in linear ICA recovery of the sources is possible up to a scalar multiplicative ambiguity.

## B WHY DOES CLASSIFICATION RESULT IN THE LOG RATIO?

Let us suppose that a variable  $X$  is drawn with equal probability from two distributions  $P_0$  and  $P_1$  with densities  $p_0(x)$  and  $p_1(x)$  respectively. We train a classifier  $D : x \mapsto [0, 1]$  to estimate the posterior probability that a particular realization of  $X$  was drawn from  $P_0$  with the cross entropy loss, i.e. the parameters of  $D$  are chosen to minimize

$$L(D) = \mathbb{E}_{X \sim P_0} [-\log D(X)] + \mathbb{E}_{X \sim P_1} [-\log(1 - D(X))].$$

As shown in, for instance, [14], the global optimum of this loss occurs when  $D(x) = \frac{p_0(x)}{p_0(x) + p_1(x)}$ , which can be rewritten as

$$D(x) = \frac{1}{1 + p_1(x)/p_0(x)} \quad (24)$$

$$= \frac{1}{1 + \exp(-\log(p_0(x)/p_1(x)))} \quad (25)$$

$$(26)$$

Recall that in our setting, the function  $r(x_1, x_2)$  is trained to classify between the two cases that  $(x_1, x_2)$  is drawn from the joint distribution  $\mathbb{P}_{x_1, x_2}$  (class 0) or the product of marginals  $\mathbb{P}_{x_1} \mathbb{P}_{x_2}$  (class 1).  $r(x_1, x_2)$  is trained so that  $\frac{1}{1 + \exp(-r(x_1, x_2))}$  estimates the posterior probability of  $(x_1, x_2)$  belonging to class 0. By comparing to Equation 25, it can be seen that

$$\begin{aligned} r(x_1, x_2) &= \log(p(x_1, x_2)/p(x_1)p(x_2)) \\ &= \log p(x_1|x_2) - \log p(x_1) \\ &= \log p(x_2|x_1) - \log p(x_2) \end{aligned}$$

Note that in order for the classification trick of contrastive learning to be useful, the variables  $x_1$  and  $x_2$  cannot be deterministically related. If this is the case, the log-ratio is everywhere either 0 or  $\infty$  and hence the learned features are not useful.

To see why this is the case, suppose that  $x_1$ , and  $x_2$  are each  $N$ -dimensional vectors. If they are deterministically related,  $p(x_1, x_2)$  puts mass on an  $N$ -dimensional submanifold of a  $2N$ -dimensional space. On the other hand,  $p(x_1)p(x_2)$  will put mass on a  $2N$ -dim manifold since it is the product of two distributions each of which are  $N$ -dimensional.

In this case, the distributions  $p(x_1, x_2)$  and  $p(x_1)p(x_2)$  are therefore not absolutely continuous with respect to one another and thus the log-ratio is ill-defined:  $p(x_1, x_2)/p(x_1)p(x_2) = \infty$  at any point  $(x_1, x_2)$  at which  $p(x_1, x_2)$  puts mass and zero at points where  $p(x_1)p(x_2)$  puts mass and  $p(x_1, x_2)$  does not.

## C THE SUFFICIENTLY DISTINCT VIEWS ASSUMPTION

We give the following two examples to provide intuition about the Sufficiently Distinct Views (SDV) assumption - one regarding a case in which it does not hold, and another one in which it does.

A simple case in which the assumption does not hold is when the conditional probability of  $z$  given  $s$  is Gaussian, as in

$$p(z|s) = \frac{1}{Z} \exp \left[ - \sum_i (z_i - s_i)^2 / (2\sigma_i^2) \right], \quad (27)$$

where  $Z$  is the normalization factor,  $Z = (2\pi)^{n/2} \prod_i \sigma_i$ . Since taking second derivatives of the log-probability with respect to  $s_i$  results in constants, it can be easily shown that there is no way to find  $2D$  vectors  $z_j$ ,  $j = 1, \dots, 2D$ , such that the corresponding  $w(s, z_j)$  (see Definition 1) are linearly independent.

The fact that the assumption breaks down in this case is reminiscent of the breakdown in the case of Gaussianity for linear ICA. Interestingly, in our work, the true latent

sources **are** allowed to be Gaussian. In fact, the distribution of  $\mathbf{s}$  does not enter the expression above.

An example in which the SDV assumption does hold is a conditional pdf given by

$$p(\mathbf{z}|\mathbf{s}) = \frac{1}{Z(\mathbf{s})} \exp \left[ - \sum_i (z_i^2 s_i^2 + z_i^4 s_i^4) \right], \quad (28)$$

where  $Z(\mathbf{s})$  is again a normalization function. Proving that this distribution satisfies the SDV assumption requires a few lines of computation. The idea is that  $\mathbf{w}(\mathbf{s}, \mathbf{z})$  can be written as the product of a matrix and vector which are functions only of  $\mathbf{s}$  and  $\mathbf{z}$  respectively. Once written in this form, it is straightforward to show that the columns of the matrix are linearly independent for almost all values of  $\mathbf{s}$  and that  $2D$  linearly independent vectors can be realized by different choices of  $\mathbf{z}$ .

## D PROOF OF THEOREM 1 AND COROLLARY 3

### D.1 PROOF OF THEOREM 1

This proof is mainly inspired by the techniques employed by [23].

*Proof.* We have to show that, upon convergence,  $h_i(\mathbf{x}_1)$  are s.t.

$$h_i(\mathbf{x}_1) \perp h_j(\mathbf{x}_1), \forall i \neq j$$

We start by writing the difference in log-densities of the two classes:

$$\begin{aligned} \sum_i \psi_i(h_i(\mathbf{x}_1), \mathbf{x}_2) &= \sum_i \alpha_i(\mathbf{f}_{1,i}^{-1}(\mathbf{x}_1), \mathbf{f}_{2,i}^{-1}(\mathbf{x}_2)) + \\ &\quad - \sum_i \delta_i(\mathbf{f}_{2,i}^{-1}(\mathbf{x}_2)) \end{aligned}$$

We now make the change of variables

$$\begin{aligned} \mathbf{y} &= \mathbf{h}(\mathbf{x}_1) \\ \mathbf{v}(\mathbf{y}) &= \mathbf{f}_1^{-1}(\mathbf{h}^{-1}(\mathbf{y})) \\ \mathbf{t} &= \mathbf{f}_2^{-1}(\mathbf{x}_2) \end{aligned}$$

and rewrite the first equation in the following form:

$$\sum_i \psi_i(y_i, \mathbf{x}_2) = \sum_i \alpha_i(v_i(\mathbf{y}), t_i) \quad (29)$$

$$- \sum_i \delta_i(t_i) \quad (30)$$

We take derivatives with respect to  $y_j, y_{j'}, j \neq j'$ , of the LHS and RHS of equation [38]. Adopting the conventions in [9] and [10] and

$$v_i^j(\mathbf{y}) = \partial v_i(\mathbf{y}) / \partial y_j \quad (31)$$

$$v_i^{jj'}(\mathbf{y}) = \partial^2 v_i(\mathbf{y}) / \partial y_j \partial y_{j'}, \quad (32)$$

we have

$$\begin{aligned} \sum_i \alpha_i''(v_i(\mathbf{y}), t_i) v_i^j(\mathbf{y}) v_i^{j'}(\mathbf{y}) \\ + \alpha_i'(v_i(\mathbf{y}), t_i) v_i^{jj'}(\mathbf{y}) = 0, \end{aligned}$$

where taking derivative w.r.t.  $y_j$  and  $y_{j'}$  for  $j \neq j'$  makes LHS equal to zero, since the LHS has functions which depend only one  $y_i$  each. If we now rearrange our variables by defining vectors  $\mathbf{a}_i(\mathbf{y})$  collecting all entries  $v_i^j(\mathbf{y}) v_i^{j'}(\mathbf{y}), j = 1, \dots, n, j' = 1, \dots, j-1$ , and vectors  $\mathbf{b}_i(\mathbf{y})$  with the variables  $v_i^{jj'}(\mathbf{y}), j = 1, \dots, n, j' = 1, \dots, j-1$ , the above equality can be rewritten as

$$\begin{aligned} \sum_i \alpha_i''(v_i(\mathbf{y}), t_i) \mathbf{a}_i(\mathbf{y}) \\ + \alpha_i'(v_i(\mathbf{y}), t_i) \mathbf{b}_i(\mathbf{y}) = 0. \end{aligned}$$

The above expression can be recast in matrix form,

$$\mathbf{M}(\mathbf{y}) \mathbf{w}(\mathbf{y}, \mathbf{t}) = 0,$$

where  $\mathbf{M}(\mathbf{y}) = (\mathbf{a}_1(\mathbf{y}), \dots, \mathbf{a}_n(\mathbf{y}), \mathbf{b}_1(\mathbf{y}), \dots, \mathbf{b}_n(\mathbf{y}))$  and  $\mathbf{w}(\mathbf{y}, \mathbf{t}) = (\alpha_1'', \dots, \alpha_n'', \alpha_1', \dots, \alpha_n')$ .  $\mathbf{M}(\mathbf{y})$  is therefore a  $n(n-1)/2 \times 2n$  matrix, and  $\mathbf{w}(\mathbf{y}, \mathbf{t})$  is a  $2n$  dimensional vector.

To show that  $\mathbf{M}(\mathbf{y})$  is equal to zero, we invoke the SDV assumption. This implies the existence of  $2n$  linearly independent  $\mathbf{w}(\mathbf{y}, \mathbf{t}_j)$ . It follows that

$$\mathbf{M}(\mathbf{y}) [\mathbf{w}(\mathbf{y}, \mathbf{t}_1), \dots, \mathbf{w}(\mathbf{y}, \mathbf{t}_{2n})] = 0,$$

and hence  $\mathbf{M}(\mathbf{y})$  is zero by elementary linear algebraic results. It follows that  $v_i^j(\mathbf{y}) \neq 0$  for at most one value of  $j$ , since otherwise the product of two non-zero terms would appear in one of the entries of  $\mathbf{M}(\mathbf{y})$ , thus rendering it non-zero. Thus  $v_i$  is a function only of one  $y_j$ .

Observe that  $\mathbf{v}(\mathbf{y}) = \mathbf{s}$ . We have just proven that  $v_i(y_{\pi(i)}) = s_i$ . Since  $v_i$  is invertible, it follows that  $h_{\pi(i)}(\mathbf{x}_1) = y_{\pi(i)} = v_i^{-1}(s_i)$  and hence the components of  $\mathbf{h}(\mathbf{x}_1)$  recover the components of  $\mathbf{s}$  up to the invertible component-wise ambiguity given by  $\mathbf{v}$ , and the permutation ambiguity. □

## D.2 PROOF OF COROLLARY 3

*Proof.* This follows exactly by repeating the proof of Theorem 1 where the roles of  $\mathbf{x}_1$  and  $\mathbf{x}_2$  are exchanged and the regression function in the statement of the corollary is used.  $\square$

## E PROOF OF THEOREMS 4 AND 5

Theorem 4 is a special case of Theorem 5 by considering the case  $\mathbf{g}_1(\mathbf{s}, \mathbf{n}_1) = \mathbf{s}$ . We therefore prove only the more general Theorem 5.

*Proof.* We have to show that, upon convergence,  $h_i(\mathbf{x}_1)$  and  $k_i(\mathbf{x}_2)$  are such that

$$h_{1,i}(\mathbf{x}_1) \perp h_{1,j}(\mathbf{x}_1), \forall i \neq j \quad (33)$$

$$h_{2,i}(\mathbf{x}_2) \perp h_{2,j}(\mathbf{x}_2), \forall i \neq j \quad (34)$$

$$h_{1,i}(\mathbf{x}_1) \perp h_{2,j}(\mathbf{x}_2), \forall i \neq j. \quad (35)$$

We start by exploiting Equations 14 and 15 to write the difference in log-densities of the two classes

$$\begin{aligned} & \sum_i \psi_i(h_{1,i}(\mathbf{x}_1), h_{2,i}(\mathbf{x}_2)) \\ &= \sum_i \eta_i(\mathbf{f}_{1,i}^{-1}(\mathbf{x}_1), \mathbf{f}_{2,i}^{-1}(\mathbf{x}_2)) - \sum_i \theta_i(\mathbf{f}_{1,i}^{-1}(\mathbf{x}_1)) \end{aligned} \quad (36)$$

$$= \sum_i \lambda_i(\mathbf{f}_{2,i}^{-1}(\mathbf{x}_2), \mathbf{f}_{1,i}^{-1}(\mathbf{x}_1)) - \sum_i \mu_i(\mathbf{f}_{2,i}^{-1}(\mathbf{x}_2)) \quad (37)$$

We now make the change of variables

$$\begin{aligned} \mathbf{y} &= \mathbf{h}_1(\mathbf{x}_1) \\ \mathbf{t} &= \mathbf{h}_2(\mathbf{x}_2) \\ \mathbf{v}(\mathbf{y}) &= \mathbf{f}_1^{-1}(\mathbf{h}_1^{-1}(\mathbf{y})) \\ \mathbf{u}(\mathbf{t}) &= \mathbf{f}_2^{-1}(\mathbf{h}_2^{-1}(\mathbf{t})) \end{aligned}$$

and rewrite equation 36 in the following form:

$$\begin{aligned} & \sum_i \psi_i(y_i, t_i) \\ &= \sum_i \eta_i(v_i(\mathbf{y}), u_i(\mathbf{t})) - \sum_i \theta_i(v_i(\mathbf{y})) \end{aligned} \quad (38)$$

We first want to prove the condition in Equation 33. We will show this is true by proving that

$$v_i(\mathbf{y}) \equiv v_i(y_{\pi(i)}) \quad (39)$$

for some permutation of the indices  $\pi$  with respect to the indexing of the sources  $\mathbf{s} = (s_1, \dots, s_D)$ .

We take derivatives with respect to  $y_j, y_{j'}, j \neq j'$ , of the LHS and RHS of equation 38, yielding

$$\begin{aligned} & \sum_i \eta_i''(v_i(\mathbf{y}), u_i(\mathbf{t})) v_i^j(\mathbf{y}) v_i^{j'}(\mathbf{y}) \\ & + \sum_i \eta_i'(v_i(\mathbf{y}), u_i(\mathbf{t})) v_i^{jj'}(\mathbf{y}) = 0 \end{aligned}$$

If we now rearrange our variables by defining vectors  $\mathbf{a}_i(\mathbf{y})$  collecting all entries  $v_i^j(\mathbf{y}) v_i^{j'}(\mathbf{y}), j = 1, \dots, n, j' = 1, \dots, j-1$ , and vectors  $\mathbf{b}_i(\mathbf{y})$  with the variables  $v_i^j(\mathbf{y}) v_i^{j'}(\mathbf{y}), j = 1, \dots, n, j' = 1, \dots, j-1$ , the above equality can be rewritten as

$$\begin{aligned} & \sum_i \eta_i''(v_i(\mathbf{y}), u_i(\mathbf{t})) \mathbf{a}_i(\mathbf{y}) \\ & + \eta_i'(v_i(\mathbf{y}), u_i(\mathbf{t})) \mathbf{b}_i(\mathbf{y}) = 0. \end{aligned}$$

Again following [23], we recast the above formula in matrix form,

$$\mathbf{M}(\mathbf{y}) \mathbf{w}(\mathbf{y}, \mathbf{t}) = 0, \quad (40)$$

where  $\mathbf{M}(\mathbf{y}) = (\mathbf{a}_1(\mathbf{y}), \dots, \mathbf{a}_n(\mathbf{y}), \mathbf{b}_1(\mathbf{y}), \dots, \mathbf{b}_n(\mathbf{y}))$  and  $\mathbf{w}(\mathbf{y}, \mathbf{t}) = (\eta_1'', \dots, \eta_n'', \eta_1', \dots, \eta_n')$ .  $\mathbf{M}(\mathbf{y})$  is therefore a  $n(n-1)/2 \times 2n$  matrix, and  $\mathbf{w}(\mathbf{y}, \mathbf{t})$  is a  $2n$  dimensional vector.

To show that  $\mathbf{M}(\mathbf{y})$  is equal to zero, we invoke the SDV assumption on  $\eta$ . This implies the existence of  $2n$  linearly independent  $\mathbf{w}(\mathbf{y}, \mathbf{t}_j)$ . It follows that

$$\mathbf{M}(\mathbf{y}) [\mathbf{w}(\mathbf{y}, \mathbf{t}_1), \dots, \mathbf{w}(\mathbf{y}, \mathbf{t}_{2n})] = 0,$$

and hence  $\mathbf{M}(\mathbf{y})$  is zero by elementary linear algebraic results. It follows that  $v_i^j(\mathbf{y}) \neq 0$  for at most one value of  $j$ , since otherwise the product of two non-zero terms would appear in one of the entries of  $\mathbf{M}(\mathbf{y})$ , thus rendering it non-zero. Thus  $v_i$  is a function only of one  $y_j = y_{\pi(i)}$ .

Observe that  $\mathbf{v}(\mathbf{y}) = \mathbf{s}$ . We have just proven that  $v_i(y_{\pi(i)}) = s_i$ . Since  $v_i$  is invertible, it follows that  $h_{\pi(i)}(\mathbf{x}_1) = y_{\pi(i)} = v_i^{-1}(s_i)$  and hence the components of  $\mathbf{h}(\mathbf{x}_1)$  recover the components of  $\mathbf{s}$  up to the invertible component-wise ambiguity given by  $\mathbf{v}$ , and the permutation ambiguity.

For the condition in Equation 34, we need

$$u_i(\mathbf{t}) \equiv u_i(t_{\tilde{\pi}(i)}), \quad (41)$$

where the permutation  $\tilde{\pi}$  doesn't need to be equal to  $\pi$ . By symmetry, exactly the same argument as used

to prove the condition in Equation 39 holds, by replacing  $(\mathbf{v}, \mathbf{y}, \boldsymbol{\eta}, \boldsymbol{\theta})$  with  $(\mathbf{u}, \mathbf{t}, \boldsymbol{\lambda}, \boldsymbol{\mu})$ , noting that the SDV assumption is also assumed for  $\boldsymbol{\lambda}$ .

We have shown that  $\mathbf{y} = \mathbf{h}_1(\mathbf{x}_1)$  and  $\mathbf{t} = \mathbf{h}_2(\mathbf{x}_2)$  estimate  $\mathbf{g}_1(\mathbf{s}, \mathbf{n}_1)$  and  $\mathbf{g}_2(\mathbf{s}, \mathbf{n}_2)$  up to two different gauges of all possible scalar invertible functions.

A remaining ambiguity could be that the two representations might be misaligned; that is, defining  $\mathbf{z}_1 = \mathbf{g}_1(\mathbf{s}, \mathbf{n}_1)$  and  $\mathbf{z}_2 = \mathbf{g}_2(\mathbf{s}, \mathbf{n}_2)$ , while

$$z_{1,i} \perp\!\!\!\perp z_{2,j} \forall i \neq j \quad (42)$$

we might have

$$y_{\pi(i)} \perp\!\!\!\perp t_{\tilde{\pi}(j)} \forall i \neq j,$$

where  $\pi(i), \tilde{\pi}(i)$  are two different permutations of the indices  $i = 1, \dots, n$ . We want to show that this ambiguity is also resolved; that means, our goal is to show that

$$y_i \perp\!\!\!\perp t_j, \quad \forall i \neq j \quad (43)$$

We recall that, by definition, we have  $v_i(y_{\pi(i)}) = z_{1,i}$  and  $u_j(t_{\tilde{\pi}(j)}) = z_{2,j}$ . Then, due to equation 42,

$$v_i(y_{\pi(i)}) \perp\!\!\!\perp u_j(t_{\tilde{\pi}(j)}) \quad \forall i \neq j \quad (44)$$

$$\implies y_{\pi(i)} \perp\!\!\!\perp t_{\tilde{\pi}(j)} \quad \forall i \neq j \quad (45)$$

$$\implies y_i \perp\!\!\!\perp t_{\tilde{\pi} \circ \pi^{-1}(j)} \quad \forall i \neq j, \quad (46)$$

where the implication 44,45 follows from invertibility of  $v_i$  and  $u_j$ , and the implication 45,46 follows from considering that, given that we know 45, we can define  $l = \pi(j)$  and  $k = \pi(i)$  and have

$$y_k \perp\!\!\!\perp t_{\tilde{\pi} \circ \pi^{-1}(l)} \quad \forall k \neq l.$$

Define

$$\tau = \tilde{\pi} \circ \pi^{-1}$$

and note that it is a permutation. Then

$$y_i \perp\!\!\!\perp t_{\tau(j)} \forall i \neq j \quad (47)$$

Fix any particular  $i$ . Our goal is to show that for any  $j \neq i$  the independence relation in Equation 43 holds. There are two possibilities:

i  $\tau(i) = i$

ii  $\tau(i) \neq i$

In the first case,  $\tau$  restricted to the set  $\{1, \dots, D\} \setminus \{i\}$  is still a permutation, and thus considering the independences of Equation 47 for all  $j \neq i$  implies each of the independences of Equation 43 and we are done.

Let us consider the second case. Then,

$$\exists l \in \{1, \dots, D\} \setminus \{i\} \text{ s.t. } l = \tau(i).$$

We then need to prove

$$y_i \perp\!\!\!\perp t_l, \quad (48)$$

which is the only independence implied by Equation 43 which is not implied by Equation 47

In order to do so, we rewrite equation 38, yielding

$$\begin{aligned} & \sum_m \psi_m(y_m, t_m) \\ &= \sum_m \eta_m(v_m(y_{\pi(m)}), u_m(t_{\tilde{\pi}(m)})) - \sum_m \theta_i(v_m(y_{\pi(m)})) \end{aligned} \quad (49)$$

We now take derivative with respect to  $y_i$  and  $t_l$  in 48; noting that  $\tilde{\pi}^{-1}(l) = \pi^{-1}(i)$ , we get

$$\begin{aligned} 0 &= \frac{\partial^2}{\partial v_{\pi^{-1}(i)} \partial u_{\pi^{-1}(i)}} \eta_{\pi^{-1}(i)}(v_{\pi^{-1}(i)}(y_i), u_{\pi^{-1}(i)}(t_l)) \\ &\quad \times \frac{\partial}{\partial y_i} v_{\pi^{-1}(i)}(y_i) \frac{\partial}{\partial t_l} u_{\pi^{-1}(i)}(t_l) \end{aligned} \quad (50)$$

Since  $v_{\pi^{-1}(i)}(y_i)$  is a smooth and invertible function of its argument, the set of  $y_i$  such that  $\frac{\partial}{\partial y_i} v_{\pi^{-1}(i)}(y_i) = 0$  has measure zero. Similarly,  $\frac{\partial}{\partial t_l} u_{\pi^{-1}(i)}(t_l) = 0$  on a set of measure zero.

It therefore follows that

$$\frac{\partial}{\partial y_i} v_{\pi^{-1}(i)}(y_i) \frac{\partial}{\partial t_l} u_{\pi^{-1}(i)}(t_l) \neq 0$$

almost everywhere and hence that

$$\frac{\partial^2}{\partial v_{\pi^{-1}(i)} \partial u_{\pi^{-1}(i)}} \eta_{\pi^{-1}(i)}(v_{\pi^{-1}(i)}(y_i), u_{\pi^{-1}(i)}(t_l)) = 0. \quad (51)$$

almost everywhere. We can thus conclude that

$$\begin{aligned} & \eta_{\pi^{-1}(i)}(v_{\pi^{-1}(i)}(y_i), u_{\pi^{-1}(i)}(t_l)) = \\ & \eta_{\pi^{-1}(i)}^y(v_{\pi^{-1}(i)}(y_i)) + \eta_{\pi^{-1}(i)}^t(u_{\pi^{-1}(i)}(t_l)) \end{aligned}$$

This in turn implies that, for some functions  $A$  and  $B$ , we can write

$$\begin{aligned} & \log p(z_{1,\pi^{-1}(i)} | z_{2,\pi^{-1}(i)}) - \log p(z_{1,\pi^{-1}(i)}) \\ &= A(v_{\pi^{-1}(i)}(y_i)) + B(u_{\pi^{-1}(i)}(t_l)) \end{aligned}$$

and therefore

$$\log p(z_{1,\pi^{-1}(i)}, z_{2,\pi^{-1}(i)}) = C(v_{\pi^{-1}(i)}(y_i)) + D(u_{\pi^{-1}(i)}(t_l))$$

for some functions  $C$  and  $D$ . This decomposition of the log-pdf implies

$$\begin{aligned} z_{1,\pi^{-1}(i)} &\perp\!\!\!\perp z_{2,\pi^{-1}(i)} \\ \implies z_{1,\pi^{-1}(i)} &\perp\!\!\!\perp z_{2,\tilde{\pi}^{-1}(l)} \\ \implies v_{\pi^{-1}(i)}(y_i) &\perp\!\!\!\perp u_{\tilde{\pi}^{-1}(l)}(t_l) \\ \implies y_i &\perp\!\!\!\perp t_l, \end{aligned}$$

where the last implication holds due to invertibility of  $v_{\pi^{-1}(i)}$  and  $u_{\tilde{\pi}^{-1}(l)}$ .

We have thus concluded the proof.  $\square$

## F PROOF OF COROLLARY 6

*Proof.* Denoting by  $d_1^{(k)}$  the component-wise invertible ambiguity up to which  $g(s, \mathbf{n}_1^{(k)})$  is recovered, we have that

$$\begin{aligned} &\inf_{e \in \mathbf{E}} \mathbb{E}_{\mathbf{x}_1} \left[ \left\| s - e(\mathbf{h}_1^{(k)}(\mathbf{x}_1)) \right\|_2^2 \right] \quad (52) \\ &= \inf_{e \in \mathbf{E}} \mathbb{E}_{(\mathbf{n}_1^{(k)}, s)} \left[ \left\| s - e \circ d_1^{(k)} \circ g_1(s, \mathbf{n}_1^{(k)}) \right\|_2^2 \right] \quad (53) \end{aligned}$$

$$= \inf_{\tilde{e} \in \mathbf{E}} \mathbb{E}_{(\mathbf{n}_1^{(k)}, s)} \left[ \left\| s - \tilde{e} \circ g_1(s, \mathbf{n}_1^{(k)}) \right\|_2^2 \right] \quad (54)$$

$$\leq \mathbb{E}_{(\mathbf{n}_1^{(k)}, s)} \left[ \left\| s - e^* \circ g_1(s, \mathbf{n}_1^{(k)}) \right\|_2^2 \right] \quad (55)$$

The lower bound holds for any  $e^* \in \mathbf{E}$  by definition of infimum and in particular for  $e^* = g_1|_{\mathbf{n}=0}^{-1}$ , the existence of which is guaranteed by the assumptions on  $g_1$ . Taking a Taylor expansion of  $e^* \circ g_1(s, \mathbf{n}_1^{(k)})$  around  $\mathbf{n}_1^{(k)} = 0$  yields

$$\begin{aligned} &\mathbb{E}_{(\mathbf{n}_1^{(k)}, s)} \left[ \left\| s - e^* \circ g_1(s, 0) \right. \right. \\ &\quad \left. \left. + \frac{\partial e^*}{\partial g_1} \frac{\partial g_1(s, 0)}{\partial \mathbf{n}_1^{(k)}} \cdot \mathbf{n}_1^{(k)} + \mathcal{O}(\|\mathbf{n}_1^{(k)}\|^2) \right\|_2^2 \right] \\ &= \mathbb{E}_{(\mathbf{n}_1^{(k)}, s)} \left[ \left\| \frac{\partial e^*}{\partial g_1} \frac{\partial g_1(s, 0)}{\partial \mathbf{n}_1^{(k)}} \cdot \mathbf{n}_1^{(k)} + \mathcal{O}(\|\mathbf{n}_1^{(k)}\|^2) \right\|_2^2 \right] \\ &\longrightarrow 0 \text{ as } k \longrightarrow \infty \end{aligned}$$

where the last equality follows from fact that  $e^* = g_1|_{\mathbf{n}=0}^{-1}$  and the convergence follows from the fact that  $\mathbf{n}_1^{(k)} \longrightarrow 0$  as  $k \rightarrow \infty$ .  $\square$

## G PROOF OF LEMMA 7

We will make crucial use of *Kolmogorov's strong law*:

**Theorem 9.** *Suppose that  $X_n$  is a sequence of independent (but not necessarily identically distributed) random variables with*

$$\sum_{n=1}^{\infty} \frac{1}{n^2} \text{Var}[X_n] < \infty$$

Then,

$$\frac{1}{N} \sum_{n=1}^N X_n - \mathbb{E}[X_n] \xrightarrow{a.s.} 0$$

Fix  $s$  and consider  $\Omega_e^N(s, \mathbf{n})$  as a random variable with randomness induced by  $\mathbf{n}$ . We will show that for almost all  $s$  this converges  $\mathbf{n}$ -almost surely to a constant, and hence  $\Omega_e^N(s, \mathbf{n})$  converges almost surely to a function of  $s$ .

The law of total expectation says that

$$\begin{aligned} &\text{Var}_{s, \mathbf{n}_i} [e_i \circ \mathbf{k}_i(s + \mathbf{n}_i)] \\ &= \mathbb{E}_s [V_i(s)] + \text{Var}_s [\mathbb{E}_{\mathbf{n}_i} [e_i \circ \mathbf{k}_i(s + \mathbf{n}_i)]] \\ &\geq \mathbb{E}_s [V_i(s)]. \end{aligned}$$

Since by assumption  $\text{Var}_{s, \mathbf{n}_i} [e_i \circ \mathbf{k}_i(s + \mathbf{n}_i)] \leq K$ , we have that

$$\mathbb{E}_s \left[ \sum_{i=1}^{\infty} \frac{V_i(s)}{i^2} \right] \leq \frac{K\pi^2}{6}$$

and therefore  $\sum_{i=1}^{\infty} \frac{V_i(s)}{i^2} < \infty$  with probability 1 over  $s$ , else the expectation above would be unbounded since  $V_i(s) \geq 0$ .

We have further that for almost all  $s$ ,

$$\Omega_e(s) = \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{i=1}^N E_{e_i}(s)$$

exists. Therefore, for almost all  $s$  the conditions of Kolmogorov's strong law are met by  $\Omega_e^N(s, \mathbf{n})$  and so

$$\Omega_e^N(s, \mathbf{n}) - \mathbb{E}_{\mathbf{n}} [\Omega_e^N(s, \mathbf{n})] \xrightarrow{n \rightarrow \infty, a.s.} 0$$

Since  $\mathbb{E}_{\mathbf{n}} [\Omega_e^N(s, \mathbf{n})] \xrightarrow{n \rightarrow \infty, a.s.} \Omega_e(s)$ , it follows that

$$\Omega_e^N(s, \mathbf{n}) \xrightarrow{n \rightarrow \infty, a.s.} \Omega_e(s).$$

Since this holds with probability 1 over  $s$ , we have that

$$\Omega_e^N(s, \mathbf{n}) \xrightarrow{n \rightarrow \infty, a.s.} \Omega_e(s).$$

It follows that we can write

$$\begin{aligned} R_{e,i}^N(s, \mathbf{n}) &= e_i \circ \mathbf{k}_i(s + \mathbf{n}_i) - \Omega_e^N(s, \mathbf{n}) \\ &\xrightarrow{a.s.} R_{e,i}(s, \mathbf{n}_i) := e_i \circ \mathbf{k}_i(s + \mathbf{n}_i) - \Omega_e(s) \end{aligned}$$

## H PROOF OF THEOREM 8

We will begin by showing that if  $K \geq \text{Var}(s) + C$  then  $\{\mathbf{k}_i^{-1}\} \in \mathcal{G}_K$ .

For  $e_i = \mathbf{k}_i^{-1}$ , we have that

$$\begin{aligned}\Omega_e^N(s, \mathbf{n}) &= \frac{1}{N} \sum_{i=1}^N s + \mathbf{n}_i \xrightarrow{a.s.} s = \Omega_e^N(s) \\ R_i^N &= s + \mathbf{n}_i - \Omega_e(s, \mathbf{n}) \xrightarrow{a.s.} \mathbf{n}_i = R_{e,i}(\mathbf{n}_i)\end{aligned}$$

where the convergences follow from application of Kolmogorov's strong law, using the fact that  $\text{Var}(\mathbf{n}_i) \leq C$  for all  $i$ . Satisfaction of condition 17 follows from the fact that  $\text{Var}_{s, \mathbf{n}_i}(s + \mathbf{n}_i) \leq C + \text{Var}(s) \leq K$ . Since  $s$  is a well-defined random variable,  $\Omega_e(s) < \infty$  with probability 1, satisfying condition 18. It follows from the mutual independence of  $\mathbf{n}_i$  and  $\mathbf{n}_j$  that  $R_{e,i}$  and  $R_{e,j}$  satisfy condition 19. Condition 20 follows from the fact that  $\mathbb{E}[\mathbf{n}_i] = 0$ . Condition 21 follows from  $R_{e,i}$  being constant as a function of  $s$ .

It therefore follows that  $\{\mathbf{k}_i^{-1}\} \in \mathcal{G}_K$  for  $K$  sufficiently large.

We will next show that if  $\{e_i\} \in \mathcal{G}_K$  then there exist a matrix  $\alpha$  and vector  $\beta$  such that  $e_i = \alpha \mathbf{k}_i^{-1} + \beta$  for all  $i$ . Since  $e_i$  acts coordinate-wise, it moreover follows that  $\alpha$  is diagonal.

First, we will show that each  $e_i \circ \mathbf{k}_i$  is affine, i.e. there exist potentially different  $\alpha_i, \beta_i$  such that  $e_i = \alpha_i \mathbf{k}_i^{-1} + \beta_i$  for each  $i$ .

Then we will show that we must have  $\alpha_i = \alpha_j$  and  $\beta_i = \beta_j$  for all  $i, j$ .

To see that  $e_i$  is affine, we make use of that fact that  $R_{e,i}$  is constant as a function of  $s$ . It follows that for any  $x$  and  $y$

$$\begin{aligned}e_i \circ \mathbf{k}_i(x + y) &= R_{e,i}(x) + \Omega_e(y) \\ &= R_{e,i}(x) + \Omega_e(0) + R_{e,i}(0) + \Omega_e(y) \\ &\quad - (R_{e,i}(0) + \Omega_e(0)) \\ &= e_i \circ \mathbf{k}_i(x) + e_i \circ \mathbf{k}_i(y) - e_i \circ \mathbf{k}_i(0)\end{aligned}$$

It therefore follows that  $e_i \circ \mathbf{k}_i$  is affine, since if we define

$$\begin{aligned}L(x + y) &= e_i \circ \mathbf{k}_i(x + y) - e_i \circ \mathbf{k}_i(0) \\ &= (e_i \circ \mathbf{k}_i(x) - e_i \circ \mathbf{k}_i(0)) \\ &\quad + (e_i \circ \mathbf{k}_i(y) - e_i \circ \mathbf{k}_i(0)) \\ &= L(x) + L(y)\end{aligned}$$

then  $L$  is linear and we can write  $e_i \circ \mathbf{k}_i(x)$  as the sum of a linear function and a constant:

$$e_i \circ \mathbf{k}_i(x) = L(x) + e_i \circ \mathbf{k}_i(0)$$

Thus  $e_i \circ \mathbf{k}_i$  is affine, and we have some (diagonal) matrix  $\alpha_i$  and vector  $\beta_i$  such that for any  $x$

$$\begin{aligned}e_i \circ \mathbf{k}_i(x) &= \alpha_i x + \beta_i \\ \implies e_i(x) &= \alpha_i \mathbf{k}_i^{-1} x + \beta_i.\end{aligned}$$

Next we show that for the set of  $\{e_i = \alpha_i \mathbf{k}_i^{-1} + \beta_i\}$ , it must be the case that each  $\alpha_i = \alpha_j$  and  $\beta_i = \beta_j$ .

Observe that

$$\begin{aligned}\Omega_e^N(s, \mathbf{n}) &= \frac{1}{N} \sum_{i=1}^N \alpha_i s + \alpha_i \mathbf{n}_i + \beta_i \\ &= \left( \frac{1}{N} \sum_{i=1}^N \alpha_i \right) s + \frac{1}{N} \sum_{i=1}^N \beta_i + \frac{1}{N} \sum_{i=1}^N \alpha_i \mathbf{n}_i \\ \mathbb{E}_{\mathbf{n}}[\Omega_e^N(s, \mathbf{n})] &= \left( \frac{1}{N} \sum_{i=1}^N \alpha_i \right) s + \frac{1}{N} \sum_{i=1}^N \beta_i\end{aligned}$$

Define

$$\begin{aligned}\alpha &= \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{i=1}^N \alpha_i \\ \beta &= \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{i=1}^N \beta_i\end{aligned}$$

which exist by the assumption that  $\Omega_e^N(s, \mathbf{n})$  converges as  $N \rightarrow \infty$ . Thus

$$\begin{aligned}\Omega_e(s) &= \alpha s + \beta \\ R_{e,i}(s, \mathbf{n}_i) &= (\alpha_i - \alpha) s + \alpha_i \mathbf{n}_i + \beta_i - \beta\end{aligned}$$

Now, suppose that there exist  $i$  and  $j$  such that  $\alpha_i \neq \alpha_j$ . It follows that

$$\begin{aligned}R_{e,i}(s, \mathbf{n}_i) &= (\alpha_i - \alpha) s + \alpha_i \mathbf{n}_i + \beta_i - \beta \\ R_{e,j}(s, \mathbf{n}_j) &= (\alpha_j - \alpha) s + \alpha_j \mathbf{n}_j + \beta_j - \beta\end{aligned}$$

There are two cases. If  $\alpha_i \neq \alpha$ , then  $R_{e,i}(s, \mathbf{n}_i)$  is not a constant function of  $s$ . But if  $\alpha_i = \alpha$ , then  $\alpha_j \neq \alpha$  and so  $R_{e,j}(s, \mathbf{n}_j)$  is not a constant function of  $s$ . This is a contradiction, and so  $\alpha_i = \alpha_j$  for all  $i, j$ .

Suppose similarly that there exist  $\beta_i \neq \beta_j$ . If  $\beta_i \neq \beta$ , then  $\mathbb{E}[R_{e,i}(\mathbf{n}_i)] = \beta_i - \beta$  which is non-zero. If  $\beta_i = \beta$ , then  $\beta_j \neq \beta$  and so  $\mathbb{E}[R_{e,j}(\mathbf{n}_j)] = \beta_j - \beta$  is non-zero. This is a contradiction, and so  $\beta_i = \beta_j$  for all  $i, j$ .

We have thus proven that set  $\{e_i\} \in \mathcal{G}_K$  is of the form  $e_i = \alpha \mathbf{k}_i^{-1} + \beta$  for all  $i$ .