

Broader impacts The ability to identify the true latent generative model for data has both positive and negative implications. On the one hand, it could be used to interpret data generative processes which can be scrutinized for to reduce implicit bias and improve algorithmic fairness. On the other hand, it may have negative consequences if privacy is a concern. Environmental concerns related to training large generative models is also worth considering.

Outline This supplement is organized as follows:

- Appendix **A** provides additional comparisons to related work, examples, and counterexamples.
- Appendix **B** provides an outline of the overall proof.
- Appendices **C-F** contain the main technical proofs.
- Appendix **G** compares the notion of equivalence in iVAE to affine equivalence.
- Appendix **H** provides conditions that guarantee a ReLU NN satisfies (F3) and/or (F4).
- Appendix **J** provides details on the experiments.

Appendices **G-H** are not required for the main proofs, but are included for completeness and reference.

A Additional details from Section 3

In this appendix, we provide additional detailed comparisons against previous work as well as additional examples and counterexamples to further illustrate our theory.

A.1 Detailed comparisons

Since the original iVAE paper [35], there have been many generalizations and extensions proposed. We pause here to provide a more detailed comparison of our results against this developing literature. For a comparison against iVAE, see Section 2.

We first discuss related work that assumes auxiliary information is available (i.e. U is known), then discuss more recent work that does not assume any auxiliary information; the ensuing comparisons are then presented in alphabetical order.

Assuming auxiliary information is available.

1. [22] achieves identifiability in the fully unsupervised regime for the model in which the latent state is defined by a Hidden Markov Model (HMM). The proof of identifiability in [22] invokes [19] to essentially recover the HMM transition matrix and the auxiliary variable U from X , reducing the problem to [35]. Our Theorem C.1 shows that identifiability in fully unsupervised regime is possible even without additional structure given here by the time-dependency according to Markov dynamics.
2. [36] extend [35] by observing that the conditional independence $Z_i \perp\!\!\!\perp Z_j \mid U$ is not required for identifiability, so they propose a more general IMCA framework for conditional energy-based models. However, identifiability in [36] still critically relies on observing an auxiliary variable (in their setting, this is a dependent variable Y). Our Theorem C.1 achieves same type of identifiability as [36] (up to affine transformation) without relying on conditional independence or an auxiliary variable.
3. [61] extends the iVAE identifiability theory of [35] by showing that a stronger notion of identifiability can be achieved if Z is distributed according to factorial GMM (instead of a general exponential family as in [35]). More specifically, given the auxiliary information U , they show that Z can be recovered up to permutation and scaling of the variables Z_i . By contrast, in Theorem E.4, we show that under similar assumptions Z_i are identifiable up to permutation and scaling and importantly, we do this only from X , without using U in any way. We also do not require the GMM to be factorial. Finally, our proof technique is different: While [61] relies on [35] (and hence, for instance, require $J \geq m + 1$), our proof is independent of [35].

4. [73] studies identifiability of the model (3) under the assumption that f is volume preserving and Z comes from a conditionally factorial exponential family, similar to iVAE. They prove that if U is known, (P2) holds, and f is twice differentiable, then Z is identifiable up to permutation and non-linear functions applied to each Z_i (i.e., $Z_i = h_i(Z_{\tau(i)})$). If additionally Z is a GMM, then Z can be recovered up to permutation, scaling, and translation. In comparison, we do not require U to be known, and we do not require f to be volume preserving or even differentiable everywhere. We show that under the same assumption (P2) the latent variables Z_i can be recovered up to permutation, scaling, and translation if f is only assumed to be piecewise affine. Additionally, we show that a weaker notion of identifiability holds if Z is not assumed to be conditionally factorial.
5. [75] considers a contrastive model in which samples arrive in pairs, which is a type of weak supervision. Additionally, it is assumed that the latent variables are sampled uniformly from a convex body, and that f is differentiable and injective. By comparison, our model allows for more general non-uniform mixture priors, non-injective and non-smooth f , and is fully unsupervised.

No auxiliary information.

1. [18] propose a novel Multifacet VAE (MFCVAE) model for unsupervised deep clustering. Their model has the following form

$$p(x, z, u) = p(x|z) \prod_{j=1}^k p(z_j|u_j)p(u_j), \quad z_j|u_j \sim \mathcal{N}(\mu_{u_j}, \Sigma_{u_j}) \quad (5)$$

Through empirical experiments, [18] emphasizes the importance of high-dimensional structure of U and shows how it results in improved clustering performance. The key idea is that while the number of meaningful clusters in the data may be very large, there may be meaningful individual categorical variables U_i (“facets”) with a much smaller number of states, which may be easier to learn. In this way, by simultaneously performing clustering for each “facet” U_i one can learn meaningful fine-grained clusters in the data. Note that k binary variables U_i result in $J = 2^k$ fine-grained clusters in the data.

Compared to our work, [18] is focused on practical implementation details, and lacks a formal identifiability theory. In fact, our results provide precisely such a formal identifiability theory in a more general setting. If $p(x|z)$ is modeled by ReLU/leaky-ReLU NN, MFCVAE is a special case of our model (3) with high-dimensional U . More specifically, the MFCVAE model (5) restricts our model (3) to the case when u_i are independent and $\text{ne}(U_i) = \{i\}$. In particular, it satisfies assumption (P3). Therefore, Theorem 3.10 implies that for MFCVAE with diagonal covariances Σ_{u_j} , $\dim(U)$, $\dim(U_j)$, $P(U)$ are identifiable from $P(X)$ up to a permutation of U , and $P(Z)$ is identifiable up to permutation, scaling, and/or translation.

2. [38] establishes the identifiability of latent representations for non-parametric measurement models $U \rightarrow X$. Their result crucially relies on the fact that observed variables are conditionally independent $X_i \perp\!\!\!\perp X_j \mid U$. Our Theorem F.2 significantly generalizes this result, by showing the same guarantees for the model (3) that allows arbitrarily complex dependencies between the observed variables X .
3. [51] propose a sparse VAE and prove that the latent space of this model is identifiable. Similar to [69], identifiability of f is not addressed. Their identifiability results also assume an anchor feature assumption, which we do not require. Even our strongest assumption (P3) is weaker compared to the anchor feature assumption (see Remark 3.8). Moreover, we do not require any sparsity assumptions.
4. [69] propose LIDVAE as a way to identify the latent space of a VAE without auxiliary information, however, their approach only guarantees identifiability of $P(Z)$, and does not address f (this is acknowledged by the authors in their discussion as an open question). By restricting f to be a Brenier map, they guarantee that the likelihood is injective, which leads to identifiability of $P(Z)$. Compared to [69] our work restricts f in a different way (i.e. by an injective ReLU network), which matches common practice. Moreover, we show that both f and the multivariate U structure (i.e. in addition to $P(Z)$) are identifiable under mild additional assumptions.

A.2 Special cases

Our main results contain some notable special cases that warrant additional discussion.

Classical VAE The classical, vanilla VAE [37, 56] with an isotropic Gaussian prior is equivalent to (3) with $J = 1$. In this case, U is trivial and the Gaussian distribution $P(Z)$ can be transformed by an affine map to a standard isotropic Gaussian $\mathcal{N}(0, I)$. In this case, Theorem 3.9(c) shows that f is identifiable from $P(X)$ up to an orthogonal transformation. In fact, this case can readily be deduced from known results on the identifiability of ReLU networks, e.g. [62].

Although the $J = 1$ case is already identifiable, there are clear reasons to prefer a clustered latent space: It is natural to model data that has several clusters by a latent space that has similar clusters (e.g. Figure 2). Although in principle any distribution can be approximated by $f(Z)$ where $Z \sim \mathcal{N}(0, I)$ and f is piecewise affine, such f is likely to be extremely complex. At the same time, the same distribution may have a representation with Z being a simple GMM and f being a simple piecewise affine function. Clearly, the latter representation is preferable to the former and can likely be more robustly learned in practice. This is consistent with previous empirical work [16, 18, 32, 33, 41, 44, 71].

Linear ICA In classical linear ICA [11], we observe $X = AZ$, where Z is assumed to have independent components. Compared to the general model (1), this corresponds to the special case where f is linear and $\varepsilon = 0$. In our most general setting under (F2) only, our results imply that $P(Z)$ can be recovered up to an affine transformation *without* assuming independent components, which might seem surprising at first. This is, however, easily explained: In this case, X is also a GMM, and hence $P(Z)$ can already be trivially recovered up to the affine transformation $z \mapsto Az$. This follows from well-known identifiability results for GMMs [63]. This provides some intuition to how the mixture prior assumption (P1) helps to achieve identifiability.

Nonlinear ICA In classical nonlinear ICA, one assumes the model (1) with (a) no assumptions on f and (b) independence assumptions in the latent space. It is well-known that this model is nonidentifiable [28]. Our problem setting is distinguished from the classical nonlinear ICA model via assumptions (P1)-(F1). While we do not require the Z_i to be mutually independent, we impose assumptions on the form of f . It is precisely this inductive bias that allows us to recover identifiability. As a result, our identifiability theory does not contradict known results such as the Darmois construction [15] discussed in [28].

A.3 Counterexamples

A natural question is whether or not the mixture prior (P1) or the piecewise affine nonlinearity (F1) can be relaxed while still maintaining identifiability. In fact, it is not hard to show this is not possible: If either (P1) or (F1) is broken, then the model (1) becomes nonidentifiable. Of course, this is entirely expected given known negative results on nonlinear ICA [28].

Example A.1. If f is allowed to be arbitrary, but (P1) is still enforced, then (1) is no longer identifiable: Pick any two GMMs $P = \sum_{j=1}^J \lambda_j N(\mu_j, \Sigma_j)$ and $P' = \sum_{j=1}^{J'} \lambda'_j N(\mu'_j, \Sigma'_j)$. Then we can always find a function g such that $g_\# P' = f_\# P$ (e.g. use the inverse CDF transform), and $g \neq f$.

Example A.2. If $P(Z)$ is allowed to be arbitrary, but (F1) is still enforced, then (1) is no longer identifiable: Consider any two arbitrary piecewise affine, injective functions $f, g : \mathbb{R}^m \rightarrow \mathbb{R}^m$. Then almost surely the preimages $f^{-1}(\{x\})$ and $g^{-1}(\{x\})$ will not be equivalent up to an affine transformation. In other words, fixing $P(X)$, we can find models (f, P) and (g, P') such that $f_\# P = P(X) = g_\# P'$, but f is not equivalent to g (i.e. up to any affine transformation).

B Proof outline

We will prove the main results by breaking the argument into four phases:

1. (Appendix C) First, we show that if f is weakly injective, then $\mathbb{P}(Z)$ is identifiable (Theorem C.1). The proof involves a novel result on identifiability of a nonparametric mixture model (Theorem C.2) that may be of independent interest.

2. (Appendix D) Second, we show that if f is continuous and injective, then f is identifiable up to an affine transformation (Theorem D.1). This result strengthens existing identifiability results in nonlinear ICA by exploiting the mixture prior, which is crucial in the sequel.
3. (Appendix E) Next, we show that if Z is conditionally factorial GMM, then under mild generic assumptions, the individual variables Z_i can be recovered (up to permutation, scaling and translation) (Theorem E.1).
4. (Appendix F) Finally, since for conditionally factorial Z we are able to recover the individual variables Z_i , we show how we can apply the theory developed in [38] to recover the multivariate discrete latent variable U , its dimension, domain sizes of each U_i and $\Pr(U, Z)$ (Theorem F.2). Since we can only recover Z up to permutation, scaling and translation, the results from [38] cannot be applied directly, and we show how to perform this recovery under an unknown affine transformation.

Each of these phases tackles a particular level of the identifiability hierarchy described in the main theorems. A detailed proof outline of each main theorem is provided below; technical proofs can be found in the subsequent appendices.

A notable difference between Theorems 3.9 ($k = 1$) and 3.10 ($k > 1$) is the conclusion in the latent space: Theorem 3.9 identifies $P(U, Z)$ jointly whereas Theorem 3.10 identifies $P(Z)$ and $P(U)$ separately. The reason is simple: If U is 1-dimensional, i.e., $k = 1$, then $P(U, Z)$ for (3) is trivially identifiable from $P(Z)$, since $P(Z)$ is assumed to be a GMM by (P1). Indeed, since finite mixture of Gaussians are identifiable, we can recover $P(U = u)$ and $P(Z | U = u)$ as mixture weights and corresponding Gaussian components. This extends to more general exponential mixtures as in Remark 2.1, see Barndorff-Nielsen [8] for details.

When $k > 1$, the situation is considerably more nontrivial, as one also needs to learn the high-dimensional structure of U .

Proof of Theorem 3.9. We assume $\varepsilon = 0$ without any loss of generality; i.e. it is sufficient to consider the noiseless case. This follows from a standard deconvolution argument as in Khemakhem et al. [36] (see Step I of the proof of Theorem 1).

- (a) By Theorem C.1, $\mathbb{P}(Z)$ is identifiable up to an affine transformation. Moreover, as described above, we can identify $\mathbb{P}(U, Z)$ from $\mathbb{P}(Z)$.
- (b) Since $P(Z)$ is identifiable up to an affine transformation by part a), claim follows from Theorem E.1.
- (c) By Theorem D.1, f is identifiable. □

Proof of Theorem 3.10. As with Theorem 3.9, we assume $\varepsilon = 0$ without loss of generality.

- (a) By Theorem C.1, $\mathbb{P}(Z)$ is identifiable up to an affine transformation.
- (b) Since $P(Z)$ is identifiable up to an affine transformation by part a), by Theorem E.1, Z_i are identifiable up to permutation, scaling and translation.
- (c) Follows from Theorem F.2.
- (d) By Theorem D.1, f is identifiable. □

C Identifiability of Z up to an affine transformation via nonparametric mixtures

In this section we prove that if in model (3) the function f is weakly injective, then Z is identifiable up to an affine transformation. More specifically, we prove the following:

Theorem C.1. Assume that (U, Z, X) are distributed according to model (3). If f is weakly injective (see (F2) in Definition 3.2), then $\mathbb{P}(U, Z)$ is identifiable from $\mathbb{P}(X)$ up to an affine transformation.

We will prove this result by first proving a result on identifiability of nonparametric mixtures that may be of independent interest.

Theorem C.2. *Let $f, g : \mathbb{R}^m \rightarrow \mathbb{R}^n$ be piecewise affine functions satisfying (F2). Let $Y \sim \sum_{i=1}^J \lambda_i \mathcal{N}(\mu_i, \Sigma_i)$ and $Y' \sim \sum_{j=1}^{J'} \lambda'_j \mathcal{N}(\mu'_j, \Sigma'_j)$ be a pair of GMMs (in reduced form). Suppose that $f(Y)$ and $g(Y')$ are equally distributed.*

Then there exists an invertible affine transformation $h : \mathbb{R}^m \rightarrow \mathbb{R}^m$ such that $h(Y) \equiv Y'$, i.e., $J = J'$ and for some permutation $\tau \in S_J$ we have $\lambda_i = \lambda'_{\tau(i)}$ and $h_{\#} \mathcal{N}(\mu_i, \Sigma_i) = \mathcal{N}(\mu'_{\tau(i)}, \Sigma'_{\tau(i)})$.

In other words, a mixture model whose components are piecewise affine transformations of a Gaussian is identifiable. To see this more clearly, observe that

$$\sum_{j=1}^J \lambda_j f_{\#} \mathcal{N}(\mu_j, \Sigma_j) \sim f_{\#} \left(\sum_{j=1}^J \lambda_j \mathcal{N}(\mu_j, \Sigma_j) \right).$$

To the best of our knowledge, this identifiability result for a nonparametric mixture model is new to the literature. In Theorem C.2, the transformation and number of components is allowed to be unknown and arbitrary, and no separation or independence assumptions are needed.

C.1 Technical lemmas

We recall that a m -dimensional Gaussian distribution $\mathcal{N}(\mu, \Sigma)$ with covariance Σ and mean μ has the following density function

$$p(x) = \frac{1}{\sqrt{(2\pi)^m \det \Sigma}} \exp \left((-1/2)(x - \mu)^T \Sigma^{-1} (x - \mu) \right). \quad (6)$$

We assume that all Gaussian components are *non-degenerate* in the sense that Σ is positive definite. We also recall that if $Y \sim \mathcal{N}(\mu, \Sigma)$ and $Y' = AY + b$ for an invertible $A \in \mathbb{R}^{m \times m}$ and $b \in \mathbb{R}^m$, then $Y' \sim \mathcal{N}(A\mu + b, A\Sigma A^T)$.

Definition C.3. We say that a Gaussian mixture distribution

$$P = \sum_{j=1}^J \lambda_j \mathcal{N}(\mu_j, \Sigma_j) \quad (7)$$

is in reduced form if $\lambda_j > 0$ for every $j \in [J]$ and for every $i \neq j \in [J]$ we have $(\mu_i, \Sigma_i) \neq (\mu_j, \Sigma_j)$.

In the proofs we use the notion of real analytic functions. We remind the definition for reader's convenience.

Definition C.4. Let $D \subseteq \mathbb{R}^n$ be an open set. A function $f : D \rightarrow \mathbb{R}$ is called a (real) analytic function if for every compact $K \subset D$ there exists a constant $C > 0$ such that for any $\alpha \in \mathbb{N}^n$ we have

$$\sup_{x \in K} \left| \frac{\partial^\alpha f}{\partial x^\alpha}(x) \right| \leq \alpha! C^{|\alpha|+1}. \quad (8)$$

Alternatively, a real analytic function $f : D \rightarrow \mathbb{R}$ can be defined as a function that has a Taylor expansion convergent on D .

It is a standard fact that a linear combination and a product of analytic functions are analytic, and it is well-known that the density of the multivariate Gaussian is a real analytic function on \mathbb{R}^m . We will also need the standard notion of analytic continuation:

Definition C.5. Let $D_0 \subseteq D \subseteq \mathbb{R}^n$ be open sets. Let $f_0 : D_0 \rightarrow \mathbb{R}$. We say that an analytic function $f : D \rightarrow \mathbb{R}$ is an *analytic continuation* of f_0 onto D if $f(x) = f_0(x)$ for every $x \in D_0$.

Definition C.6. Let $x_0 \in \mathbb{R}^m$ and $\delta > 0$. Let $p : B(x_0, \delta) \rightarrow \mathbb{R}$. Define

$$\text{Ext}(p) : \mathbb{R}^m \rightarrow \mathbb{R} \quad (9)$$

to be the unique analytic continuation of p on the entire space \mathbb{R}^m if such a continuation exists, and to be 0 otherwise.

Definition C.7. Let $D_0 \subset D$ and $p : D \rightarrow \mathbb{R}$ be a function. We define $p|_{D_0} : D_0 \rightarrow \mathbb{R}$ to be a restriction of p to D_0 , namely a function that satisfies $p|_{D_0}(x) = p(x)$ for every $x \in D_0$.

Theorem C.8. Consider a pair of finite GMMs (in reduced form) in \mathbb{R}^m

$$P = \sum_{j=1}^J \lambda_j \mathcal{N}(\mu_j, \Sigma_j) \quad \text{and} \quad P' = \sum_{j=1}^{J'} \lambda'_j \mathcal{N}(\mu'_j, \Sigma'_j). \quad (10)$$

Assume that there exists a ball $B(x_0, \delta)$ such that P and P' induce the same measure on $B(x_0, \delta)$. Then $P \equiv P'$, i.e., $J = J'$ and for some permutation τ we have $\lambda_i = \lambda'_{\tau(i)}$ and $(\mu_i, \Sigma_i) = (\mu'_{\tau(i)}, \Sigma'_{\tau(i)})$.

Proof. Follows from the identity theorem for real analytic functions and the identifiability of finite GMMs. \square

Definition C.9. Let $f : \mathbb{R}^m \rightarrow \mathbb{R}^n$ be a piecewise affine function. We say that a point $x \in f(\mathbb{R}^m) \subseteq \mathbb{R}^n$ is *generic with respect to f* if the preimage $f^{-1}(\{x\})$ is finite and there exists $\delta > 0$, such that $f : B(z, \delta) \rightarrow \mathbb{R}^n$ is affine for every $z \in f^{-1}(\{x\})$.

Lemma C.10. If $f : \mathbb{R}^m \rightarrow \mathbb{R}^n$ is a piecewise affine function such that $\{x \in \mathbb{R}^n : |f^{-1}(\{x\})| = \infty\} \subseteq f(\mathbb{R}^m)$ has measure zero with respect to the Lebesgue measure on $f(\mathbb{R}^m)$, then $\dim(f(\mathbb{R}^m)) = m$ and almost every point in $f(\mathbb{R}^m)$ (with respect to the Lebesgue measure on $f(\mathbb{R}^m)$) is generic with respect to f .

Proof. Let $g_i(z) = Az + b$, $g : D \rightarrow \mathbb{R}^n$ be one of the affine pieces defining piecewise affine function f . If A does not have full column rank, then every $x \in g(D)$ has an infinite number of preimages. Therefore, the assumption of the lemma implies that for at least one of the affine pieces g_i , A has full column rank. Thus, $\dim(f(\mathbb{R}^m)) = m$.

Let $S = \{x \in \mathbb{R}^n : |f^{-1}(\{x\})| = \infty\}$ then by assumption S has measure zero in $f(\mathbb{R}^m)$. Let E be the set of points $z \in \mathbb{R}^m$ such that for every $\delta > 0$, f is not affine on $B(z, \delta)$. Since f is piecewise affine, E can be covered by a locally-finite union of $(m-1)$ -dimensional subspaces, i.e. every compact set intersects only finitely many of these (potentially infinite) $(m-1)$ -dimensional subspaces. Thus E has measure zero. Moreover, since $\dim(f(\mathbb{R}^m)) = m$, $f(E)$ has measure zero in $f(\mathbb{R}^m)$.

Finally, by definition, every $x \in f(\mathbb{R}^m) \setminus (S \cup f(E))$ is generic. \square

We make the following useful observation.

Lemma C.11. Consider a random variable Z distributed according to the GMM $\sum_{j=1}^J \lambda_j \mathcal{N}(\mu_j, \Sigma_j)$.

Consider the random variable $X = f(Z)$, where $f : \mathbb{R}^m \rightarrow \mathbb{R}^m$ is a piecewise affine function, such that $\dim(f(\mathbb{R}^m)) = m$. Let $x_0 \in \mathbb{R}^m$ be a generic point with respect to f . Let p be the density function of X . Then the number of points in the preimage $f^{-1}(\{x_0\})$ can be computed as

$$|f^{-1}(\{x_0\})| = \lim_{\delta \rightarrow 0} \int_{x \in \mathbb{R}^m} \text{Ext}(p|_{B(x_0, \delta)})(x) dx. \quad (11)$$

Proof. Since x_0 is generic with respect to f , the preimage of x_0 consists of finitely many points, $f^{-1}(\{x_0\}) = \{z_1, z_2, \dots, z_s\}$, and there exists $\varepsilon > 0$ such that for every $i \in [s]$ there is a well-defined invertible affine function $g_i : B(z_i, \varepsilon) \rightarrow \mathbb{R}^m$ such that $g_i(z) = f(z)$ for all $z \in B(z_i, \varepsilon)$.

We can write $g_i(z) = A_i z + b_i$ for some $A_i \in \mathbb{R}^{m \times m}$ and $b_i \in \mathbb{R}^m$. Let $\delta_0 > 0$ be such that

$$B(x_0, \delta_0) \subseteq \bigcap_{i=1}^s g_i(B(z_i, \varepsilon)). \quad (12)$$

Let $0 < \delta < \delta_0$. Then, for $\mu'_{ij} = A_i \mu_j + b_i$ and $\Sigma_{ij} = A_i \Sigma_j A_i^T$, and every $x \in B(x_0, \delta)$ we have

$$p|_{B(x_0, \delta)}(x) = \sum_{i=1}^s \sum_{j=1}^J \frac{\lambda_j}{\sqrt{(2\pi)^m \det \Sigma_{ij}}} \exp\left((-1/2)(x - \mu'_{ij})^T \Sigma_{ij}^{-1} (x - \mu'_{ij})\right). \quad (13)$$

The RHS of (13) is a real analytic function defined on all of \mathbb{R}^m (i.e. it is an entire function) that equals p on an open neighborhood, hence it defines $\text{Ext}(p|_{B(x_0, \delta)})$ on the entire space \mathbb{R}^m . Therefore,

$$\begin{aligned} & \int_{x \in \mathbb{R}^m} \text{Ext}(p|_{B(x_0, \delta)})(x) dx = \\ &= \int_{x \in \mathbb{R}^m} \sum_{i=1}^s \sum_{j=1}^J \frac{\lambda_j}{\sqrt{(2\pi)^m \det \Sigma_{ij}}} \exp \left((-1/2)(x - \mu'_{ij})^T \Sigma_{ij}^{-1} (x - \mu'_{ij}) \right) = \\ &= \sum_{i=1}^s \int_{x \in \mathbb{R}^m} \sum_{j=1}^J \frac{\lambda_j}{\sqrt{(2\pi)^m \det \Sigma_{ij}}} \exp \left((-1/2)(x - \mu'_{ij})^T \Sigma_{ij}^{-1} (x - \mu'_{ij}) \right) = \\ &= s = |f^{-1}(\{x_0\})|. \end{aligned} \quad \square$$

We can deduce the following corollary.

Corollary C.12. *Let $f, g : \mathbb{R}^m \rightarrow \mathbb{R}^n$ be piecewise affine functions that satisfy (F2).*

Let $Z \sim \sum_{i=1}^J \lambda_i \mathcal{N}(\mu_i, \Sigma_i)$ and $Z' \sim \sum_{j=1}^{J'} \lambda'_j \mathcal{N}(\mu'_j, \Sigma'_j)$. Suppose that $f(Z)$ and $g(Z')$ are equally distributed. Assume that for $x_0 \in \mathbb{R}^n$ and $\delta > 0$, f is invertible on $B(x_0, 2\delta) \cap f(\mathbb{R}^m)$.

Then there exists $x_1 \in B(x_0, \delta)$ and $\delta_1 > 0$ such that both f and g are invertible on $B(x_1, \delta_1) \cap f(\mathbb{R}^m)$.

Proof. Since f is piecewise affine and f is invertible on $B(x_0, 2\delta) \cap f(\mathbb{R}^m)$, then $\dim f(\mathbb{R}^m) = m$. Note that since $f(Z)$ and $g(Z')$ are equally distributed and since regular GMMs have positive density at every point, we have

$$f(\mathbb{R}^m) = \text{supp}(f(Z)) = \text{supp}(g(Z')) = g(\mathbb{R}^m).$$

Therefore, $\dim(g(\mathbb{R}^m)) = \dim(f(\mathbb{R}^m)) = m$ and, by Lemma C.10, almost every point $x \in B(x_0, \delta) \cap f(\mathbb{R}^m)$ is generic with respect to f and w.r.t to g . Let $x_1 \in B(x_0, \delta)$ be such a point. Since f is invertible on $B(x_1, \delta)$, we have that $|f^{-1}(\{x_1\})| = 1$. Since x_1 is generic with respect to f and with respect to g , by Lemma C.11, we deduce that $|g^{-1}(\{x_1\})| = 1$. Therefore, since x_1 is generic, there exists $0 < \delta_1 < \delta$ such that on $(B(x_1, \delta_1) \cap f(\mathbb{R}^m)) \subset (B(x_0, 2\delta) \cap f(\mathbb{R}^m))$ the function g is invertible. \square

C.2 Identifiability of nonparametric mixtures

First we prove our identifiability theorem under the assumption that f and g are invertible in the neighborhood of the same point.

Theorem C.13. *Let $f, g : \mathbb{R}^m \rightarrow \mathbb{R}^n$ be piecewise affine. Let $Z \sim \sum_{i=1}^J \lambda_i \mathcal{N}(\mu_i, \Sigma_i)$ and $Z' \sim \sum_{j=1}^{J'} \lambda'_j \mathcal{N}(\mu'_j, \Sigma'_j)$ be a pair of GMMs (in reduced form). Suppose that $f(Z)$ and $g(Z')$ are equally distributed.*

Assume that there exists $x_0 \in \mathbb{R}^n$ and $\delta > 0$ such that f and g are invertible on $B(x_0, \delta) \cap f(\mathbb{R}^m)$. Then there exists an invertible affine transformation $h : \mathbb{R}^m \rightarrow \mathbb{R}^m$ such that $h(Z) \equiv Z'$, i.e., $J = J'$ and for some permutation τ we have $\lambda_i = \lambda'_{\tau(i)}$ and $h_{\#} \mathcal{N}(\mu_i, \Sigma_i) = \mathcal{N}(\mu'_{\tau(i)}, \Sigma'_{\tau(i)})$.

Proof. Since f and g are piecewise affine and both f and g are invertible on $B(x_0, \delta) \cap f(\mathbb{R}^m)$, then $\dim f(\mathbb{R}^m) = m$ and the inverse functions are piecewise affine. Hence, moreover, there exist x_1 and $\delta_1 > 0$ with $B(x_1, \delta_1) \subseteq B(x_0, \delta)$ such that f^{-1} and g^{-1} on $B(x_1, \delta_1) \subseteq B(x_0, \delta)$ are defined by affine functions.

Let $L \subseteq \mathbb{R}^n$ be an m -dimensional affine subspace, such that $B(x_1, \delta_1) \cap f(\mathbb{R}^m) = B(x_1, \delta_1) \cap L$.

Let $h_f, h_g : \mathbb{R}^m \rightarrow L$ be a pair of invertible affine functions such that h_f^{-1} coincides with f^{-1} on $B(x_1, \delta_1) \cap L$ and h_g^{-1} coincides with g^{-1} on $B(x_1, \delta_1) \cap L$. This means that distributions $h_f(Y)$

and $h_g(Y')$ coincide on $B(x_1, \delta_1) \cap L$. Moreover, since h_f and h_g are affine transformations, then $h_f(Y)$ and $h_g(Y')$ are finite GMMs. Therefore, by Theorem C.8, $h_f(Y) \equiv h_g(Y')$. The claim of the theorem holds for $h = h_g^{-1} \circ h_f$. \square

Combining this identifiability result with results of Section C.1, we obtain the proof of our main identifiability result for non-parametric mixtures.

Proof of Theorem C.2. By Corollary C.12 there exists $x_0 \in f(\mathbb{R}^m)$ that is generic with respect to both f and g and $\delta > 0$ such that f and g are invertible on $B(x_0, \delta) \cap f(\mathbb{R}^m)$. Therefore, the result follows from Theorem C.13. \square

C.3 Proof of Theorem C.1

We give a proof by contradiction. Assume that there exists another model (U', Z', X') and a piecewise affine function g in model 3 that generates the same distribution, i.e., $\mathbb{P}(X) = \mathbb{P}(X')$.

By Corollary C.12 there exists $x_0 \in f(\mathbb{R}^m)$ that is generic with respect to both f and g and $\delta > 0$ such that f and g are invertible on $B(x_0, \delta) \cap f(\mathbb{R}^m)$. Therefore, by Theorem C.13, there exists $h : \mathbb{R}^m \rightarrow \mathbb{R}^m$ such that $Z' = h(Z)$. In other words, $P(U, Z)$ is identifiable up to an affine transformation.

D Identifiability of f

In this section we show that if f is continuous piecewise affine and injective then it is identifiable from $P(X)$ up to an affine transformation.

Theorem D.1. Assume that (U, Z, X) are distributed according to model (3). Assume that f is continuous piecewise affine and satisfies (F4) (i.e., f is injective).

Then $(\mathbb{P}(U, Z), f)$ is identifiable from $\mathbb{P}(X)$ up to an affine transformation.

Before proving this theorem, we provide an example that shows that assumption (F2) does not guarantee that f can be recovered uniquely up to an affine transformation in Theorem C.1.

Example D.2. Consider

$$Y \sim \frac{1}{2}\mathcal{N}(-2, 1) + \frac{1}{2}\mathcal{N}(2, 1) \quad (14)$$

Define a pair of piecewise affine functions (see also Figure 3)

$$f(x) = \begin{cases} x - 4, & \text{for } x \geq 4, \\ -x, & \text{for } -2 \leq x < 2, \\ x + 4, & \text{for } -4 \leq x < -2, \\ (x + 4)/5, & \text{for } x < -4. \end{cases} \quad g(x) = \begin{cases} x - 4, & \text{for } x \geq 4, \\ -x + 4, & \text{for } 2 \leq x < 4, \\ x, & \text{for } -2 \leq x < 2, \\ -x - 4, & \text{for } -4 \leq x < -2, \\ (x + 4)/5, & \text{for } x < -4. \end{cases} \quad (15)$$

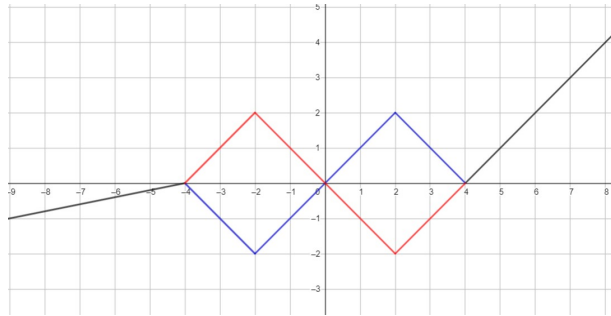


Figure 3: Graphs of f (black and red) and g (black and blue) in Example D.2.

Then it is easy to see that $f(Y)$ and $g(Y)$ have the same distribution, but f cannot be transformed into g by an affine transformation.

In order to prove Theorem D.1 we need to show that for a mixture of Gaussians P and a pair of piecewise affine functions f, g if $f_{\#}P = g_{\#}P$, then $f = h \circ g$ for some invertible affine h . We first consider the case when g is the identity.

Lemma D.3. *Let $Z \sim \sum_{j=1}^J \lambda_j \mathcal{N}(\mu_j, \Sigma_j)$. Assume that $f : \mathbb{R}^m \rightarrow \mathbb{R}^m$ is a continuous piecewise affine function such that $f(Z) \sim Z$. Then f is affine.*

Proof. Since Z has positive density at every point and $f(Z) \sim Z$ we must have $\dim f(\mathbb{R}^m) = m$.

If f is not affine, then there exist an $(m-1)$ -dimensional affine subspace L , $z_0 \in L$ and $\delta > 0$ such that the following holds: The subspace L divides $B(z_0, \delta)$ into two sets (formally, these are “half-balls”) B^+ and B^- such that $f_+(z) := f|_{B^+}(z) = A_1 z + b_1$ and $f_-(z) := f|_{B^-}(z) = A_2 z + b_2$, where $(A_1, b_1) \neq (A_2, b_2)$ and A_1, A_2 are invertible.

Since $f(Z) \sim Z$ we have

$$\{f_+(\mu_1), \dots, f_+(\mu_J)\} = \{\mu_1, \mu_2, \dots, \mu_K\} = \{f_-(\mu_1), \dots, f_-(\mu_J)\}$$

as multisets (i.e. including repetitions). Let $\mu_* = \frac{1}{J} \sum_{j=1}^J \mu_j$. Then, since f_+ and f_- are affine we get $f_+(\mu_*) = f_-(\mu_*) = \mu_*$. By translating Y and adjusting f accordingly, we may assume that $\mu_* = 0$. In this case, $b_1 = b_2 = 0$. Moreover, since $f_+(z) = f_-(z)$ for $z \in L$, we get

$$(A_1^{-1}A_2)(z) = z \quad \text{for all } z \in L. \quad (16)$$

Finally, since $f(Y) \sim Y$, we have

$$\{A_1 \Sigma_1 A_1^T, \dots, A_1 \Sigma_J A_1^T\} = \{\Sigma_1, \Sigma_2, \dots, \Sigma_J\} = \{A_2 \Sigma_1 A_2^T, \dots, A_2 \Sigma_J A_2^T\},$$

as multisets (i.e. including repetitions). This implies that

$$\prod_{j=1}^J \det(A_1 \Sigma_j A_1^T) = \prod_{j=1}^J \det(\Sigma_j) = \prod_{j=1}^J \det(A_2 \Sigma_j A_2^T).$$

Hence, $\det(A_1)^2 = \det(A_2)^2 = 1$, and $\det(A_1^{-1}A_2)^2 = 1$. By (16), $A_1^{-1}A_2$ is the identity map on L . Let v be a unit vector orthogonal to L (in the direction of B^+). Then we get that either $A_1^{-1}A_2 v = v$, or $A_1^{-1}A_2 v = -v$. In the latter case $A_1(y_0 + (\delta/2)v) = A_2(y_0 - (\delta/2)v)$, which means that f is not injective. This contradicts Lemma C.11. Therefore, we must have $A_1^{-1}A_2 v = v$, and so, by (16), $A_1 = A_2$.

Therefore, $f_+ = f_-$, which contradicts $(A_1, b_1) \neq (A_2, b_2)$. It follows that f must be affine. \square

Theorem D.4. *Let $f, g : \mathbb{R}^m \rightarrow \mathbb{R}^n$ be continuous invertible piecewise affine functions. Let $Z \sim \sum_{i=1}^J \lambda_i \mathcal{N}(\mu_i, \Sigma_i)$ and $Z' \sim \sum_{j=1}^{J'} \lambda'_j \mathcal{N}(\mu'_j, \Sigma'_j)$ be a pair of GMMs (in reduced form). Suppose that $f(Z)$ and $g(Z')$ are equally distributed.*

Then there exists an affine transformation $h : \mathbb{R}^m \rightarrow \mathbb{R}^m$ such that $h(Z) \equiv Z'$ and $g = f \circ h^{-1}$.

Proof. By Theorem C.13, there exists an invertible affine transformation $h_0 : \mathbb{R}^m \rightarrow \mathbb{R}^m$ such that $h_0(Z) = Z'$. Then, $f(Z) \sim g(h_0(Z))$, and since g and h_0 are invertible, we can rewrite this as $Z \sim (h_0^{-1} \circ g^{-1} \circ f)(Z)$. By Lemma D.3, $(h_0^{-1} \circ g^{-1} \circ f)$ is affine, i.e. there exists an invertible affine map h_1 such that

$$h_0^{-1} \circ g^{-1} \circ f = h_1 \quad \Leftrightarrow \quad f = g \circ (h_0 \circ h_1)$$

Hence the claim of the theorem holds for $h = h_0 \circ h_1$. \square

Proof of Theorem D.1. Immediately follows from Theorems C.1 and D.4. \square

D.1 Identifiability under assumption (F3)

In this section we discuss the case (F3). In particular, show that in (3) under the weaker assumption (F3), f is identifiable up to an affine transformation on the preimage of every connected open set onto which f is injective.

Theorem D.5. *Let $f, g : \mathbb{R}^m \rightarrow \mathbb{R}^n$ be continuous piecewise affine functions satisfying (F3).*

Let $Z \sim \sum_{i=1}^J \lambda_i \mathcal{N}(\mu_i, \Sigma_i)$ and $Z' \sim \sum_{j=1}^{J'} \lambda'_j \mathcal{N}(\mu'_j, \Sigma'_j)$ be a pair of variables with GMM distribution (in reduced form). Suppose that $f(Z)$ and $g(Z')$ are equally distributed.

Let $\mathcal{D} \subseteq \mathbb{R}^n$ be a connected open set such that f and g are injective onto \mathcal{D} . Then there exists an affine transformation $h : \mathbb{R}^m \rightarrow \mathbb{R}^m$ such that $h(Z) \equiv Z'$ and $g(z) = (f \circ h^{-1})(z)$ for every $z \in g^{-1}(\mathcal{D})$.

Proof. Similarly, as in the proof of Theorem D.4, by Theorem C.13, there exists an invertible affine transformation $h_0 : \mathbb{R}^m \rightarrow \mathbb{R}^m$ such that $h_0(Z) = Z'$. Then, $f(Z) \sim g(h_0(Z))$, and since g is invertible on \mathcal{D} and h_0 is invertible, we can rewrite this as $Z \sim (h_0^{-1} \circ g^{-1} \circ f)(Z)$ on $f^{-1}(\mathcal{D})$. Since f is invertible and continuous piecewise affine, $f^{-1}(\mathcal{D})$ is an open connected set. Therefore, applying Lemma D.3 on $f^{-1}(\mathcal{D})$, we deduce that $(h_0^{-1} \circ g^{-1} \circ f)$ is affine on $f^{-1}(\mathcal{D})$, i.e. there exists an invertible affine map h_1 such that

$$h_0^{-1} \circ g^{-1} \circ f = h_1 \quad \Leftrightarrow \quad f = g \circ (h_0 \circ h_1) \quad \text{on } f^{-1}(\mathcal{D})$$

Therefore, for $h = (h_0 \circ h_1)$, we have $g(y) = (f \circ h^{-1})(z)$ for every $z \in g^{-1}(\mathcal{D})$. \square

Remark D.6. Let f be a continuous piecewise affine function that satisfies (F3). Denote

$$S = \{x \in \mathbb{R}^n : |f^{-1}(\{x\})| > 1\} \subseteq f(\mathbb{R}^m).$$

Recall that assumption (F3) says that S has measure zero in $f(\mathbb{R}^m)$.

We claim that (F3) implies that for every $x \in S$ in fact $|f^{-1}(\{x\})| = \infty$. Indeed, if for all sufficiently small $\delta > 0$ we have $\dim(B(x, \delta) \cap f(\mathbb{R}^m)) < m$, then $|f^{-1}(\{x\})| = \infty$ since f is continuous piecewise affine. Otherwise, using Corollary C.10, we get that for every $\delta > 0$ there exists a generic with respect to f point $x_\delta \in B(x, \delta) \cap f(\mathbb{R}^m)$. Assumption (F3) implies that $|f^{-1}(\{x_\delta\})| = 1$ for every x_δ . Therefore, since f is continuous piecewise affine we get that either $|f^{-1}(\{x\})| = 1$ or $|f^{-1}(\{x\})| = \infty$.

E Identifiability of Z up to a permutation, scaling and translation

Under (P2), we have

$$Z \sim \sum_{j=1}^J \lambda_j \mathcal{N}(\mu_j, \Sigma_j), \quad (17)$$

where Σ_j is diagonal for every $j \in [J]$. In the setup of model (3) this just means that $Z_i \perp\!\!\!\perp Z_j \mid U$.

Let $Y = AZ + b$, where $A : \mathbb{R}^m \rightarrow \mathbb{R}^m$ is an invertible linear map and $b \in \mathbb{R}^m$. Then Y is also a GMM. We next show how Z may be recovered from Y up to a permutation, scaling, and translation.

Theorem E.1. *Let $J \geq 2$, and $\lambda_j > 0$ for all $j \in [J]$. Let $Z = (Z_1, Z_2, \dots, Z_m)$ be given by*

$$Z \sim \sum_{j=1}^J \lambda_j \mathcal{N}(\mu_j, \Sigma_j) \quad (18)$$

Assume that Σ_j is diagonal for every $j \in [J]$. Let $Y = AZ + b$, where $A : \mathbb{R}^m \rightarrow \mathbb{R}^m$ is an invertible linear map and $b \in \mathbb{R}^m$. Moreover, assume that there exist indices $i_1, i_2 \in [J]$, such that all numbers $((\Sigma_{i_1})_{tt} / (\Sigma_{i_2})_{tt} \mid t \in [m])$ are distinct. Given Y , one can recover an invertible linear map $A' : \mathbb{R}^m \rightarrow \mathbb{R}^m$, such that $(A')^{-1}A = QD$, where Q is a permutation matrix and D is a diagonal matrix with positive entries.

Remark E.2. The translation b is impossible to recover without stronger assumptions, as b corresponds to an arbitrary translation in the Z space. In other words, choice of b determines the origin in the coordinate space of Z and it can be completely arbitrary.

Remark E.3. A slightly different version of Theorem E.1 under different assumptions appeared in [73]. The main difference is that [73] assumed that f is volume-preserving but nonlinear, whereas we restrict to the general (i.e. not necessarily volume-preserving) linear case.

Proof. Without loss of generality assume $i_1 = 1$ and $i_2 = 2$.

Let Σ_i be the covariance matrices of Z_i and let $\tilde{\Sigma}_i$ be the covariance matrices of Y_i for $i \in [J]$. Clearly

$$\tilde{\Sigma}_i = A \Sigma_i A^T \quad \text{for each } i \in [J]. \quad (19)$$

The matrices $\tilde{\Sigma}_i$ are PSD. Therefore, using SVD we can find PSD matrices V_i , such that for every $i \in [J]$,

$$\tilde{\Sigma}_i = V_i V_i^T. \quad (20)$$

Moreover, such a decomposition is unique up to an orthogonal matrix, i.e., for every pair of such decompositions $\tilde{\Sigma}_i = V_i V_i^T = V'_i (V'_i)^T$ there exists a unitary matrix R such that $V_i R = V'_i$. Therefore, for every $i \in [J]$ there exists a matrix R_i , such that

$$V_i R_i = A \Sigma_i^{1/2} \quad (21)$$

In particular,

$$V_1 R_1 \Sigma_1^{-1/2} = V_2 R_2 \Sigma_2^{-1/2} \Rightarrow R_1 \left(\Sigma_1^{-1/2} \Sigma_2^{1/2} \right) R_2^{-1} = \left(V_1^{-1} V_2 \right) \quad (22)$$

Since R_1 and R_2^{-1} are unitary and $\left(\Sigma_1^{-1/2} \Sigma_2^{1/2} \right)$ is diagonal, they can be determined from the SVD of $V_1^{-1} V_2$. Moreover, they can be determined uniquely up to a permutation matrix since all diagonal entries of $\Sigma_1^{-1/2} \Sigma_2^{1/2}$ are distinct. In other words, using SVD for $\left(V_1^{-1} V_2 \right)$ we can find R'_1 such that for some permutation matrix P we have

$$V_1 R'_1 Q = A \Sigma_1^{1/2}, \quad \text{so, for } A' := V_1 R_1 \quad \text{we have } (A')^{-1} A = Q \Sigma_1^{-1/2}. \quad (23)$$

This concludes the proof. \square

As an immediate corollary we can deduce the following theorem from Theorem C.1.

Theorem E.4. Assume that (U, Z, X) are distributed according to model (3) and that f is weakly injective. Suppose that $Z_i \perp\!\!\!\perp Z_j \mid U$ for all $i \neq j$. Moreover, assume that there exist a pair of states $U = u_1$ and $U = u_2$ such that all $((\Sigma_{u_1})_{tt} / (\Sigma_{u_2})_{tt} \mid t \in [m])$ are distinct.

Then $P(U, Z)$ is identifiable from $P(X)$ up to permutation, scaling and translation of Z_i .

Proof. By Theorem C.1, $\mathbb{P}(Z)$ is identifiable from $\mathbb{P}(X)$ up to an affine transformation. That is, we can reconstruct a random variable Y from $\mathbb{P}(X)$ which satisfies $Y = AZ + b$ for some invertible $A \in \mathbb{R}^{m \times m}$.

Now, by Theorem E.1, we can find A' such that $Z' = (A')^{-1} Y = Q D Z + (A')^{-1} b$, where Q is a permutation matrix and D is a diagonal matrix. This means, that we can recover Z up to permutation, shift and scaling of individual variables Z_i . \square

F Identifiability of multivariate U structure

When $k = 1$, $P(Z)$ contains all the information about $P(U, Z)$, however, when $k > 1$ (i.e. U is multivariate), this may not be true anymore. It is not even obvious that $P(Z)$ must contain information about the true dimension of U . The distribution $P(U, Z)$ may contain interesting dependencies between individual variables U_i and Z_j .

Previously, [38] studied necessary and sufficient conditions for identifiability of $P(U)$ when Z is observed under the so-called *measurement model*. A key limitation of Kivva et al. [38] is that it

requires the observed variables to be conditionally independent, which is not the case in our setting. Ultimately, this is a consequence of Z being unobserved: Previous work such as Kivva et al. [38] assumes there is only a single layer of hidden variables connected to the observations. In our setting, under (3), we need to recover U from Z , the latter of which is unobserved. As a result, if we can only identify Z up to an affine transformation (e.g., like in Theorem C.1); i.e. we can only recover $Z' = AZ + b$, then it almost surely will not be conditionally factorial. Hence, the results from Kivva et al. [38] cannot be applied directly for weak (e.g., up to affine transformation, or as in 35) notions of identifiability of Z .

Luckily, in Section E, we showed how to recover the true Z from $Z' = AZ + b$. This will enable us to identify $P(U)$ in Theorem 3.10(c). In the remainder of this appendix, we outline these details.

We say that a distribution $\mathbb{P}(U, Z)$ satisfies the *Markov property* with respect to the neighborhoods $\text{ne}(Z_i)$ (cf. Definition 3.7) if

$$\mathbb{P}(U, Z) = \mathbb{P}(U) \prod_i \mathbb{P}(Z_i \mid \text{ne}(Z_i)). \quad (24)$$

Remark F.1. The neighborhoods $\text{ne}(Z_i)$ define a bipartite graph between (U_1, \dots, U_k) and (Z_1, \dots, Z_m) that is described in Kivva et al. [38]. Since this graph is not needed for our purposes, we proceed without further mention of this graph. The assumptions below have been re-phrased accordingly.

[38] show that assumptions (L1)-(L4) below are necessary for identifiability of U .

- (L1) (No twins) For any $U_i \neq U_j$ we have $\text{ne}(U_i) \neq \text{ne}(U_j)$.
- (L2) (Maximality) There is no U' such that:
 - (a) $\mathbb{P}(U', Z)$ is Markov with respect to the neighborhoods $\text{ne}(Z_i)$ defined by U' ;
 - (b) U' is obtained from U by splitting a hidden variable (equivalently, U is obtained from U' by merging a pair of vertices);
 - (c) U' satisfies Assumption (L1).
- (L3) (Nondegeneracy) The distribution over (U, Z) satisfies:
 - (a) $\mathbb{P}(U = u) > 0$ for all u .
 - (b) For all $Z' \subset Z$ and $u_1 \neq u_2$, $\mathbb{P}(Z' \mid \text{ne}(Z') = u_1) \neq \mathbb{P}(Z' \mid \text{ne}(Z') = u_2)$, where u_1 and u_2 are distinct configurations of $\text{ne}(Z')$.
- (L4) (Subset condition) For any pair of distinct variables U_i, U_j the set $\text{ne}(U_i)$ is not a subset of $\text{ne}(U_j)$.

We prove the following identifiability result.

Theorem F.2. Assume that (U, Z, X) are distributed as in (3) and that f satisfies (F2). Assume further that (P2)-(P3) hold and $P(U = u) > 0$ for all u in the domain of U .

Then $\dim(U) = k$, $\dim(U_j)$, $\mathbb{P}(U, Z)$ are identifiable from $P(X)$ up to a permutation of variables U_i and permutation, scaling and translation of variables Z_i .

Proof. The assumptions of Theorem F.2 are stronger than those of Theorem E.4, so by Theorem E.4, $P(Z)$ is identifiable up to a permutation, scaling and translation of Z .

Combined with the positivity assumption $P(U = u) > 0$, the assumptions (L1)-(L4) are weaker than assumption (P3). Indeed, (P3) (a) is equivalent to (L3) (b); (P3) (c) is equivalent to (L4) and implies (L1); and, finally, (P3) (b) and (c) together imply (L2).

Since Z is identifiable up to a permutation, scaling and translation, $Z_i \perp\!\!\!\perp Z_j \mid U$, and assumptions (L1)-(L4) hold, using [38, Thm 3.2], we deduce that $\dim(U) = k$, $\dim(U_j)$, $\mathbb{P}(U)$, and $\text{ne}(U_i)$ are identifiable up to a permutation of the variables U_i . Finally, by the Markov Property, $\mathbb{P}(U)$, $\text{ne}(U_i)$ for all i , and the fact that $P(Z)$ is a finite GMM (that is identifiable) are sufficient to recover $\mathbb{P}(U, Z)$. \square

Remark F.3. As the proof indicates, assumptions (L1)-(L4) are weaker than (P3), so Theorem F.2 implies part (c) of Theorem 3.10.

G Equivalence in iVAE

In this section we compare the equivalence relation up to which iVAE [35] guarantees identifiability and equivalence up to an affine transformation. While iVAE achieves the best possible identifiability under the assumptions they make, we show that identifiability up to an affine transformation is considerably stronger.

G.1 iVAE equivalence relation

Recall that iVAE [35] considers the following model, which differs from (3) by assuming that Z has conditionally factorial exponential family distribution:

$$\left. \begin{aligned} U &= u \sim p(u) \\ [Z | U = u] &\sim \prod_{i=1}^m \frac{Q_i(z_i)}{C(u)} \exp \left(\sum_{j=1}^t T_{i,j}(z_i) \lambda_{i,j}(u) \right) \\ [X | Z = z] &\sim f(z) + \varepsilon, \quad \varepsilon \sim \mathcal{N}(v, \sigma^2) \end{aligned} \right\} \implies U \rightarrow Z \rightarrow X. \quad (25)$$

Here $T_i = (T_{i,1}, T_{i,2}, \dots, T_{i,t})$ are sufficient statistics, Q_i is the base measure and $\lambda_{i,j}$ parameters depending on u . iVAE defines the following equivalence relation:

Definition G.1.

$$(f, T, \sigma) \sim (f', T', \sigma') \iff \exists A, c : T(f^{-1}(\{x\})) = A(T'((f')^{-1}(x)) + c, \quad (26)$$

where $A : \mathbb{R}^{mt} \rightarrow \mathbb{R}^{mt}$ is an invertible linear map, and $c \in \mathbb{R}^N$.

This type of identifiability allows for essentially any (synchronized) changes to Z and f :

Lemma G.2. *Let $\varphi : \mathbb{R}^m \rightarrow \mathbb{R}^m$ be any invertible map. Let $f' = f \circ \varphi$, and $T' = T \circ \varphi$. Then $(f, T, \sigma) \sim (f', T', \sigma)$.*

Moreover, if Z has exponential family distribution with statistics T , then $Z' = \varphi^{-1}(Z)$, has an exponential family distribution with statistics T' , and $f(Z) \sim f'(Z')$.

Proof. We have $(f')^{-1} = \varphi^{-1} \circ f^{-1}$, so $T' \circ (f')^{-1} = T \circ f$. Hence $(f, T, \sigma) \sim (f', T', \sigma')$, where in (26) A is the identity map and $c = 0$.

Since Z comes from an exponential family distribution, we can write

$$\mathbb{P}(Z | U) = h(Z)g(U) \exp(\lambda(U)T(Z)). \quad (27)$$

Let $Z' = \varphi^{-1}(Z)$. Then by the change of variable formula

$$\mathbb{P}(Z' | U) = \left(h(\varphi(Z)) \det |Jac(\varphi(\bullet))|_{\bullet=\varphi^{-1}(Z)} \right) g(U) \exp(\lambda(U)T(\varphi(Z))), \quad (28)$$

where $Jac(\varphi)$ is the Jacobian of φ . Hence Z' indeed has an exponential family distribution with statistics T' . Clearly, $f'(Z') = (f \circ \varphi \circ \varphi^{-1})(Z) \equiv f(Z)$. \square

Remark G.3. In other words, the equivalence relation (26) allows an *arbitrary* (possibly highly nonlinear) change of basis in the latent Z space. In principle, this may indicate, that any meaningful analysis of the Z space in this setup may be challenging.

Remark G.4. As in Khemakhem et al. [35], the additional assumption that Z has a conditionally factorial distribution imposes additional restrictions on φ . In this case, $\varphi : \mathbb{R}^m \rightarrow \mathbb{R}^m$ can be any invertible coordinatewise function $\varphi(Z') = (\varphi_1(z'_1), \varphi_2(z'_2), \dots, \varphi_m(z'_m))$.

G.2 GMMs give more robust identifiability

The next result was also observed in [61]. We present a slightly simplified proof for completeness.

If $\mathbb{P}(Z|U)$ is a multivariate Gaussian distribution, then the sufficient statistics are given by

$$T_m = (z_1, \dots, z_m, z_1 z_1, z_1 z_2, \dots, z_m z_m). \quad (29)$$

Remark G.5. For product measures, there are no cross-terms $z_i z_j$.

Proposition G.6 (61, Appendix B). Assume that $(T_m, f, \sigma) \sim (T_m, f', \sigma')$, where T_m is defined by (29). Then there exists an invertible linear map $M : \mathbb{R}^m \rightarrow \mathbb{R}^m$ and a vector $c \in \mathbb{R}^m$ such that $f^{-1}(\{x\}) = M(f')^{-1}(x) + c$ for every x .

Proof. Let $z = f^{-1}(\{x\})$ and $z' = (f')^{-1}(x)$. By an assumption of the proposition there exists an invertible matrix $A : \mathbb{R}^{m+m^2} \rightarrow \mathbb{R}^{m+m^2}$ such that

$$\begin{pmatrix} z_1 \\ z_2 \\ \vdots \\ z_n \\ z_1 z_1 \\ z_1 z_2 \\ \vdots \\ z_m z_m \end{pmatrix} = A \begin{pmatrix} z'_1 \\ z'_2 \\ \vdots \\ z'_n \\ z'_1 z'_1 \\ z'_1 z'_2 \\ \vdots \\ z'_m z'_m \end{pmatrix} + b \quad (30)$$

This means that for every i there exists a polynomial p_i of degree at most 2 such that $z_i = p_i(z'_1, \dots, z'_m)$. Assume that for some i , we have $\deg(p_i) = 2$. Then it is easy to verify (say, by using lexicographical order on monomials) that $\deg(p_i^2) = 4$. If z' is defined on an open neighbourhood, we get a contradiction with (30) as z_i^2 can be written as a degree-2 polynomial over variables z'_j . Therefore, every p_i is a polynomial of degree at most 1. But this means that $z = Mz' + c$ for some matrix M and a vector c . Moreover, since A is invertible, M is invertible as well. \square

H Conditions on ReLU Neural Network that guarantee that it is an observable injection

For completeness, in this section we provide simple sufficient conditions on ReLU architectures that guarantee that it is an observable injection (cf. (F3)) and simple sufficient conditions on leaky-ReLU architectures which guarantee that it is injection (cf. (F4)). For a more comprehensive account of identifiability in ReLU networks, see [62].

We recall the definitions of ReLU and leaky-ReLU (with parameter $a > 0$, $a \neq 1$) activation functions

$$\text{ReLU}(x) = \begin{cases} x, & \text{for } x > 0, \\ 0, & \text{for } x \leq 0, \end{cases} \quad \text{LReLU}(x) = \begin{cases} x, & \text{for } x > 0, \\ a \cdot x, & \text{for } x \leq 0. \end{cases} \quad (31)$$

A standard choice of a for leaky-ReLU is $a = 0.01$.

Definition H.1. Let $\text{Aff}(n_1, n_2)$ denote the set of affine maps $h : \mathbb{R}^{n_1} \rightarrow \mathbb{R}^{n_2}$.

Let $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ be a general activation function. For a vector $x \in \mathbb{R}^t$, $\sigma(x)$ is a vector obtained from x by applying σ coordinatewise.

Definition H.2. Let $n_1, n_2, \dots, n_t \geq n_0 = m$ and σ be an activation function. Define

$$\mathcal{F}_\sigma^{n_0, \dots, n_t} = \{h_t \circ \sigma \circ h_{t-1} \circ \sigma \circ \dots \circ h_1 \mid h_i \in \text{Aff}(n_{i-1}, n_i)\} \quad (32)$$

$$\mathcal{F}_\sigma^{m \hookrightarrow n} = \bigcup_{t=1}^{\infty} \bigcup_{n_1, n_2, \dots, n_t \geq n_0, n_0=m, n_t=n} \mathcal{F}_\sigma^{n_0, \dots, n_t} \quad (33)$$

Remark H.3. The function families $\mathcal{F}_{\text{ReLU}}^{m \hookrightarrow n}$, $\mathcal{F}_{\text{LReLU}}^{m \hookrightarrow n}$ are genuinely nonparametric: There is no bound on the number of layers.

Remark H.4. In the arguments below we do not rely on the fact that the activation function is the same on every layer, or even the same across the nodes of the same layer. However, we will give proofs only in this case, to simplify the presentation.

Remark H.5. ReLU networks under similar assumptions were also studied in [36].

Lemma H.6. *Let $f = h_t \circ \sigma \circ h_{t-1} \circ \sigma \circ \dots \circ h_1 \in \mathcal{F}_{\text{ReLU}}^{m \hookrightarrow n}$. Assume that $m = n_0 \leq n_1 \leq \dots \leq n_t = n$, and $\dim(f(\mathbb{R}^m)) = m$. Then for almost all $y \in f(\mathbb{R}^m)$ there exists δ_y such that f^{-1} is a well-defined affine function on $B(y, \delta_y) \cap f(\mathbb{R}^m)$.*

Proof. We prove the claim by induction on the depth of the NN. If $t = 1$, we have $f = h_1$ and the claim is trivial. Assume that we already proved the lemma for all $t \leq s - 1$. We prove the claim for $t = s$. We can write f as $f = h_t \circ \sigma \circ g$ where $g \in \mathcal{F}_{\text{ReLU}}^{m \hookrightarrow n_{t-1}}$.

Since $\dim(f(\mathbb{R}^m)) = m$, the map h_t has full column rank. Additionally, denoting by $\mathcal{D} = \{x \in \mathbb{R}^{n_{t-1}} \mid x_i > 0, \forall i \in [n_{t-1}]\}$ the domain on which σ is injective, we get $g(\mathbb{R}^m) \cap \mathcal{D}$ has positive measure in $g(\mathbb{R}^m)$. Moreover, by the induction assumption, g satisfies conclusion of the lemma, i.e., there exists a set S of measure 0 in $g(\mathbb{R}^m)$ such that for any $y \in g(\mathbb{R}^m) \setminus S$ there exists a $\delta_y > 0$ such that g^{-1} is a well-defined affine function on $B(y, \delta_y) \cap g(\mathbb{R}^m)$. Since h_t has full column rank, f^{-1} is a well-defined affine function on $B(x, \delta_x) \cap f(\mathbb{R}^m)$ for every $x = (f \circ \sigma)(y)$ where $y \in (g(\mathbb{R}^m) \setminus S) \cap \mathcal{D}$. Clearly, such x form a set of full measure in $f(\mathbb{R}^m)$. \square

Corollary H.7. *Let $f = h_t \circ \sigma \circ h_{t-1} \circ \sigma \circ \dots \circ h_1 \in \mathcal{F}_{\text{ReLU}}^{m \hookrightarrow n}$. Assume that $m = n_0 \leq n_1 \leq \dots \leq n_t = n$, and $\dim(f(\mathbb{R}^m)) = m$, then f satisfies (F3).*

Proof. Immediately follows from Lemma H.6. \square

Lemma H.8. *Let $f = h_t \circ \sigma \circ h_{t-1} \circ \sigma \circ \dots \circ h_1 \in \mathcal{F}_{\text{LReLU}}^{m \hookrightarrow n}$. Assume that $m = n_0 \leq n_1 \leq \dots \leq n_k = n$ and every h_i is invertible. Then for almost all $y \in f(\mathbb{R}^m)$ there exists δ_y such that f^{-1} is a well-defined affine function on $B(y, \delta_y) \cap f(\mathbb{R}^m)$.*

Proof. Clearly, any $f = h_t \circ \sigma \circ h_{t-1} \circ \sigma \circ \dots \circ h_1 \in \mathcal{F}_{\text{LReLU}}$ is a piecewise affine function. The LReLU activation function is invertible, so f is invertible. Finally, since f is a piecewise affine transformation, for almost all $y \in B(x_0, \delta)$ there exists δ_y such that f^{-1} is an affine function on $B(y, \delta_y)$. \square

Corollary H.9. *Let $f = h_t \circ \sigma \circ h_{t-1} \circ \sigma \circ \dots \circ h_1 \in \mathcal{F}_{\text{LReLU}}^{m \hookrightarrow n}$. Assume that $m = n_0 \leq n_1 \leq \dots \leq n_t = n$, then generically f satisfies (F4).*

Proof. Generically, every h_i has full column rank, and so is injective. Since LReLU is injective, we get that f is injective. \square

We conclude with an example of a very simple LReLU NN that is not even weakly injective.

Example H.10. Let $\sigma(x) = x$ for $x \geq 0$ and $\sigma(x) = x/2$ for $x < 0$. Let $h_1 : \mathbb{R} \rightarrow \mathbb{R}^2$ defined as $h_1(x) = (x, -x)$. Then $\sigma \circ h_1(x) = (x, -x/2)$ if $x \geq 0$ and $\sigma \circ h_1(x) = (x/2, -x)$ if $x < 0$. Let $h_2 : \mathbb{R}^2 \rightarrow \mathbb{R}^2$ given by

$$h_2 = \begin{pmatrix} 1 & -1 \\ 1 & 1 \end{pmatrix}$$

Then $(h_2 \circ \sigma \circ h_1)(x) = (3x/2, x/2)$ for $x \geq 0$ and $(h_2 \circ \sigma \circ h_1)(x) = (3x/2, -x/2)$ for $x < 0$ (see Figure 4). Let $h_3(x, y) = y$. Then $f(x) := (h_3 \circ \sigma \circ h_2 \circ \sigma \circ h_1)(x) = |x|/2$. By Remark 3.12, this implies that f is not invertible at every point except 0.

J Experiment details

J.1 Metrics

Previous work has relied on the Mean Correlation Coefficient (MCC) as a metric to quantify identifiability. For consistency with previous work, we report this metric, but also propose a new metric to quantify identifiability up to an affine transformation. There are two challenges in designing such a metric: Firstly, for two Gaussian mixtures, standard distance metrics such as TV-distance or KL-divergence do not have a closed form. Secondly, we need to find an affine map A that best aligns a pair of Gaussian mixtures. Therefore, developing a metric to quantify identifiability up to an affine transformation has natural challenges. We propose $\text{dist}_{\text{Aff}, L_2}$, defined below, as an additional metric in this setting.

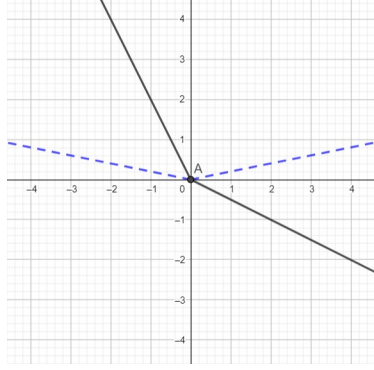


Figure 4: Graphs of $\sigma \circ h_1$ (black) and $h_2 \circ \sigma \circ h_1$ (blue) in Example H.10

Measuring loss In this work, we consider two different metrics. For a pair of distributions p_1, p_2 , we define $\text{dist}_{\text{Aff}, L_2}$ loss as

$$\text{dist}_{\text{Aff}, L_2}(p_1, p_2) = \min_{\substack{A: \mathbb{R}^m \rightarrow \mathbb{R}^m \\ \text{affine}}} \Delta_{L_2}(A_{\#} p_1, p_2), \quad \text{where} \quad \Delta_{L_2}(p_1, p_2) = \frac{\|p_1 - p_2\|_{L_2}}{\|p_1\|_{L_2}^{1/2} \|p_2\|_{L_2}^{1/2}} \quad (34)$$

The other metric we consider is the Mean Correlation Coefficient (MCC) metric which had been used in prior works [36, 70]. See Khemakhem et al. [36, Appendix A.2] for a detailed discussion. There are two versions of MCC that have been used:

- The *strong* MCC is defined to be the MCC before alignment via the affine map A .
- The *weak* MCC is defined to be the MCC after alignment.

In our experiments, we report both the strong MCC and weak MCC. Moreover, all reported MCCs are out-of-sample, i.e. the optimal affine map A is computed over half the dataset and then reused for the other half of the dataset.

Alignment To find the affine map A that best aligns the two GMMs, we use two approaches. One approach is to use Canonical Correlation Analysis (CCA) as was done in prior works in computing MCC.

We describe an alternative approach now. Given two GMMs, we iterate over all permutations of the components and for each fixed permutation, we find the best map A that maps the components accordingly. In an ideal setting, we would want to find A to align not just the means but also the covariance matrices but unfortunately this is a challenging optimization problem. Therefore, we instead find A that maps the means of the first GMM to the means of the second GMM. The map A can be found by solving a least-squares optimization problem which is straightforward using a Singular Value Decomposition (SVD). In practice, we find that this technique of matching the means works well.

J.2 Implementation

For VaDE [32], we use the implementation available at <https://github.com/mperezcarrasco/Pytorch-VaDE>. For MFCVAE [18], we use the author implementation available at <https://github.com/FabianFalck/mfcvae>. For iVAE [35], we use the implementation available at <https://github.com/MatthewWillettts/algostability>. Experiments were performed on an NVIDIA Tesla K80 GPU with 12GB memory.

J.3 Setup

Our experiments consist of three different setups, designed to probe different aspects of identifiability. First, we checked the exact log-likelihood for a unique global minimizer on simple toy models

(Appendix J.3.1). We then used VaDE [32] to train a practical VAE on a simulated dataset where the ground truth latent space is known (Appendix J.3.2). Finally, we compared the performance of MFCVAE [18] against iVAE on MNIST (Appendix J.3.3). The last experiment is based on previous work by [70] that compares iVAE to VaDE; we successfully replicated these experiments using MFCVAE as an additional baseline that closely aligns with our assumptions.

The fact that our theory closely aligns with and replicates existing empirical work illustrates that the model (3) is not merely a theoretical curiosity, but in fact practically relevant in modern applications. In our view, this is a significant advantage compared to related work.

J.3.1 Maximum likelihood

We simulated random models of the form (1) as follows:

1. Fix $J = 2$ or $J = 3$;
2. Randomly select $(\lambda_1, \dots, \lambda_J)$ from a uniform grid by discretizing the simplex;
3. Randomly select (μ_1, \dots, μ_J) from a uniform grid on the hypercube;
4. Randomly select coefficients (α_1, α_2) , weights (β_1, β_2) , and biases (π_1, π_2) from a uniform grid on the hypercube.

Given these parameters, the prior $P(Z)$ is defined as in (2) and the decoder f is defined to be the following single-layer ReLU network

$$f(z) = \alpha_1 \text{ReLU}(\beta_1 z + \pi_1) + \alpha_2 \text{ReLU}(\beta_2 z + \pi_2).$$

As a result of the simulation mechanism, the following important cases of misspecification naturally arise:

- We allow $\lambda_j = 0$, i.e. the model allows for $J = 3$ components, but the true model only has two nontrivial components.
- We allow $\alpha_j = 0$ and $\beta_j = 0$, i.e. the model allows for up to two neurons in the hidden layer, but the true model only has one nontrivial neuron.
- f is not forced to be injective or even weakly injective, i.e. assumptions (F2)-(F4) are not checked explicitly.

After generating a pair $(f, P(Z))$, the exact negative log-likelihood is approximated via numerical integration. An exhaustive grid search is performed over all parameters to identify the global minimizers. The computational cost of this step limited the complexity of the models that could be tested, hence the restriction to simple toy models in this experiment. In all runs, the ground truth was the unique global minimizer of the negative log-likelihood, as predicted by our theory. Since the problem is nonconvex, there often exist additional (non-global) local minima (see e.g. Figure 1), however, the global minimizer is always unique up to affine equivalence. That is, due to affine equivalence, in some cases there is more than one global minimizer, but in all such cases it is easy to check that the different minimizers are indeed affinely equivalent. Multiple minimizers also arise when certain parameters (e.g. λ_j or α_j) vanish, again, these are easily checked.

J.3.2 Simulated data

We consider 4 synthetic datasets described below: Pinwheel and three different copies of the “Random parallelograms” dataset

See Section 4 for results of the simulated experiments on the “pinwheels” dataset (see 33). In those experiments we use 5000 samples and set $m = n = 2$. In that experiment we used the same neural network architecture as discussed below for “Random parallelograms”.

We simulate an artificial dataset “Random parallelograms” as follows: We generate 3 randomly oriented parallelograms in the plane. After that, an n -dimensional observed distribution is obtained by sampling points uniformly at random from these parallelograms and by adding Gaussian noise to every sampled point.

We fit VaDE to each (observed) dataset 5 times (see Figures 2, 5-7). Let $Z^{(1)}, Z^{(2)}, \dots, Z^{(5)}$ be the learned latent spaces. For every pair $Z^{(i)}, Z^{(j)}$ we evaluate the MCC and $\text{dist}_{\text{Aff}, L2}$ loss. We report means of the MCCs/losses and their standard deviations in Table 2.

For the VaDE training, we use a sequential neural network architecture with LeakyReLU activations for the encoder, with four fully connected layers of the following dimensions: $n \rightarrow 64 \rightarrow 512 \rightarrow 64 \rightarrow m$. For the decoder, we use a sequential neural network architecture with LeakyReLU activations, with four fully connected layers of the following dimensions: $m \rightarrow 64 \rightarrow 512 \rightarrow 512 \rightarrow n$. We pretrain the autoencoder for 15 epochs and then run VaDE training for 20 epochs.

In all experiments with simulated data we set $m = 2$. We set the number of observed samples to be 5000.

Dataset	$\text{dist}_{\text{Aff}, L2}$	Strong MCC	Weak MCC
Random parallelograms #1	0.1542 (0.150)	0.86 (0.09)	0.99 (0.003)
Random parallelograms #2	0.1231 (0.076)	0.83 (0.12)	0.99 (0.003)
Random parallelograms #3	0.578 (0.301)	0.91 (0.08)	0.99 (0.001)

Table 2: Mean (std) $\text{dist}_{\text{Aff}, L2}$ distance (lower is better) and Mean (std) MCC (higher is better) for synthetic data

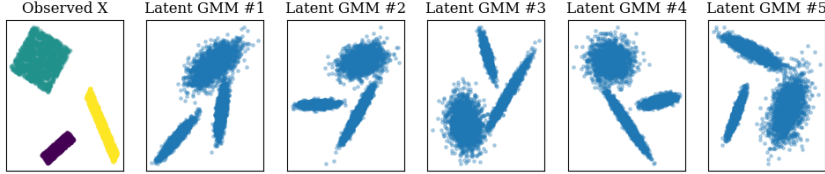


Figure 5: Recovered latent spaces for 5 runs of VaDE on “Random parallelograms” dataset #1 with 3 clusters

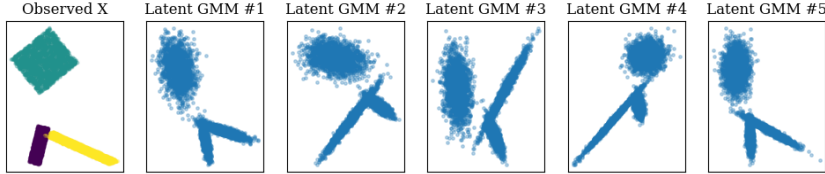


Figure 6: Recovered latent spaces for 5 runs of VaDE on “Random parallelograms” dataset #2 with 3 clusters

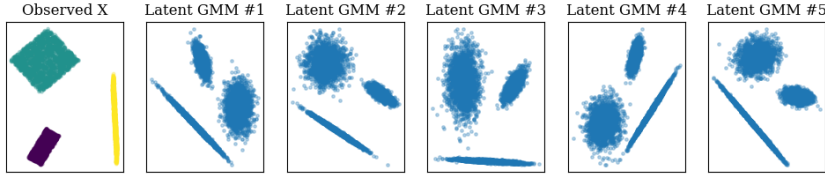


Figure 7: Recovered latent spaces for 5 runs of VaDE on “Random parallelograms” dataset #3 with 3 clusters

J.3.3 Real data

We run MFCVAE [18] on the MNIST dataset 10 times with different initializations. For all the 45 pairs of runs, we compute the strong MCC (before alignment) and weak MCC (after alignment with CCA of dimension 5). For these experiments, we omit the $\text{dist}_{\text{Aff}, L2}$ metric since it’s computationally infeasible with a large number of components. The mean and standard deviation of the MCCs are

reported in Table 3. As a baseline, we also report the same metrics for 10 runs of iVAE [35] on identical architecture and latent dimension, but recall that iVAE has additional access to the true digit labels U .

Architecture	Model	Activation	Strong MCC	Weak MCC
Arch1	MFCVAE	ReLU	0.7 (0.07)	0.91 (0.05)
	MFCVAE	LeakyReLU	0.69 (0.06)	0.94 (0.02)
	iVAE	LeakyReLU	0.65 (0.07)	0.88 (0.07)
Arch2	MFCVAE	ReLU	0.69 (0.07)	0.89 (0.08)
	MFCVAE	LeakyReLU	0.69 (0.06)	0.92 (0.03)
	iVAE	LeakyReLU	0.64 (0.07)	0.87 (0.04)
Arch3	MFCVAE	ReLU	0.69 (0.07)	0.86 (0.08)
	MFCVAE	LeakyReLU	0.70 (0.05)	0.92 (0.03)
	iVAE	LeakyReLU	0.67 (0.06)	0.87 (0.05)

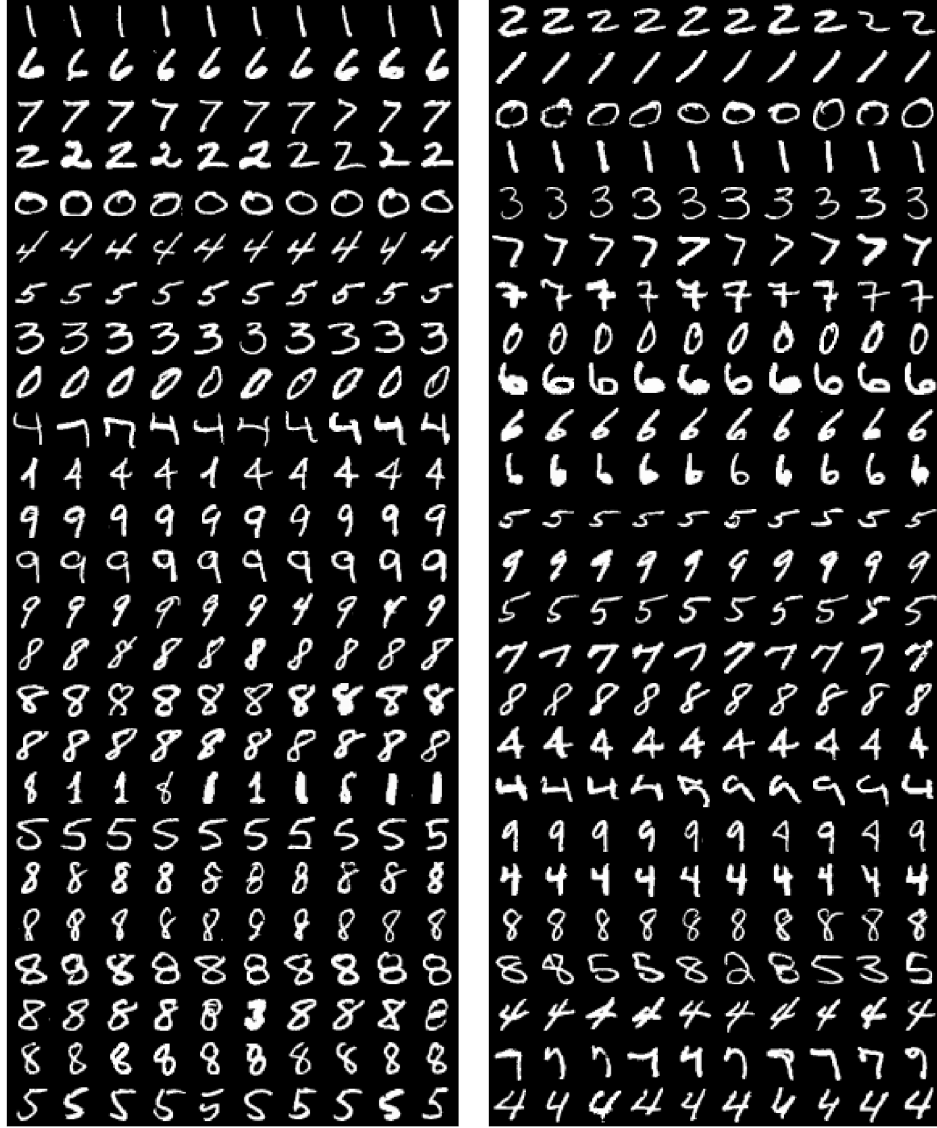
Table 3: Mean and standard deviation of the MCCs (higher is better) across various models, architectures and activations

As recommended in [18], we set the dimension of the latent space to be 5 and number of components to be 25. No hyperparameter tuning was done. The architectures we use are as follows:

- Arch1: The encoder is a sequential neural network architecture with fully connected layers of dimensions $n \rightarrow 500 \rightarrow 1000 \rightarrow m$. The decoder is also a sequential neural network architecture with fully connected layers of dimensions $m \rightarrow 500 \rightarrow 500 \rightarrow n$.
- Arch2: The encoder is a sequential neural network architecture that is fully connected with dimensions $n \rightarrow 256 \rightarrow 512 \rightarrow 512 \rightarrow m$. The decoder is similarly a sequential neural network architecture with fully connected layers of dimensions $m \rightarrow 512 \rightarrow 256 \rightarrow n$.
- Arch3: The encoder is a sequential neural network architecture that is fully connected with dimensions $n \rightarrow 128 \rightarrow 256 \rightarrow 128 \rightarrow 128 \rightarrow m$. The decoder is again a sequential neural network architecture with fully connected layers of dimensions $m \rightarrow 128 \rightarrow 128 \rightarrow n$.

The work [70] ran extensive experiments comparing VaDE and iVAE. We augment these experiments by using MFCVAE instead of VaDE. We observe that even without access to U , MFCVAE has competitive performance (stability) in recovering the latent space as compared to iVAE which has full access to U . This offers strong evidence for stability of training, as predicted by our theory.

For purely illustrative purposes, we also show the output of MFCVAE on MNIST. In Figure 8, we show samples synthetically generated from each learnt cluster. In Figure 9, we visualize the true datapoint x and the corresponding reconstructed \hat{x} for four different datapoints in each cluster. For similar experiments on other datasets and other architectures, we refer the reader to [18].



(a) Arch1

(b) Arch2

Figure 8: Output of MFCVAE on MNIST data: Synthetically generated samples. Each row corresponds to a different learnt component. The columns are samples generated from the component. The rows are sorted by average confidence.



(a) Arch1



(b) Arch2

Figure 9: Output of MFCVAE on MNIST data: Reconstruction accuracy. Each row corresponds to a different learnt component, the columns correspond to 4 different pairs of x and \hat{x} in that order.