

Supplementary Material

A Computing the Latent Space Partition

In this section we first introduce notations and demonstrate how to express a region ω of the partition Ω as a polytope defined by a system of inequalities, and then leverage this formulation to demonstrate how to obtain Ω by recursively exploring neighboring regions starting from a random point/region.

Regions as Polytopes To represent the regions $\omega \in \Omega$ as a polytope via a system of inequalities we need to recall from (1) that the input-output mapping is defined on each region by the affine parameters A_ω, B_ω themselves obtained by composition of MASOs. Each layer pre-activation (feature map prior application of the nonlinearity) is denoted by $\mathbf{h}^\ell(\mathbf{z}) \in \mathbb{R}^{D^\ell}, \ell = 1, \dots, L-1$ and given by $\mathbf{h}^\ell(\mathbf{x}) = \mathbf{A}_\omega^{1 \rightarrow \ell} \mathbf{z} + \mathbf{b}_\omega^{1 \rightarrow \ell}$, with up-to-layer ℓ affine parameters

$$\mathbf{A}_\omega^{1 \rightarrow \ell} \triangleq \mathbf{W}^\ell \mathbf{D}_\omega^{\ell-1} \mathbf{W}^{\ell-1} \dots \mathbf{D}_\omega^1 \mathbf{W}^1, \quad \mathbf{A}_\omega^{1 \rightarrow \ell} \in \mathbb{R}^{D^\ell \times S}, \quad (10)$$

$$\mathbf{b}_\omega^{1 \rightarrow \ell-1} \triangleq \mathbf{v}^\ell + \sum_{i=1}^{\ell} \mathbf{W}^\ell \mathbf{D}_\omega^{\ell-1} \mathbf{W}^{\ell-1} \dots \mathbf{D}_\omega^i \mathbf{v}^i, \quad \mathbf{b}_\omega^{1 \rightarrow \ell} \in \mathbb{R}^{D^\ell}, \quad (11)$$

which depend on the region ω in the latent space². Notice that we have in particular $\mathbf{A}_\omega^L = \mathbf{A}_\omega$ and $\mathbf{b}_\omega^L = \mathbf{b}_\omega$, the entire DGN affine parameters from (2) on region ω . The regions depend on the signs of the pre-activations defined as $\mathbf{q}^\ell(\mathbf{z}) = \text{sign}(\mathbf{h}^\ell(\mathbf{z}))$ due to the used activation function behaving linearly as long as the feature maps preserve the same sign. This holds for (leaky-)ReLU or absolute value, for max-pooling we would need to look at the argmax position of each pooling window, as pooling is rare in DGN we focus here on DN without max-pooling; let $\mathbf{q}^{\text{all}}(\mathbf{z}) \triangleq [(\mathbf{q}^{L-1}(\mathbf{z}))^T, \dots, (\mathbf{q}^1(\mathbf{z}))^T]^T$ collect all the per layer sign operators without the last layer as it does not apply any activation.

Lemma 5. *The \mathbf{q}^{all} operator is piecewise constant and there is a bijection between Ω and $\text{Im}(\mathbf{q})$.*

The above demonstrates the equivalence of knowing ω in which an input \mathbf{z} belongs to and knowing the sign pattern of the feature maps associated to \mathbf{z} ; we will thus use interchangeably $\mathbf{q}^{\text{all}}(\mathbf{z}), \mathbf{z} \in \omega$ and $\mathbf{q}^{\text{all}}(\omega)$. From this, we see that the pre-activation signs and the regions are tied together. We can now leverage this result and provide the explicit region ω as a polytope via its system of inequality, to do so we need to collect the per-layer slopes and biases into

$$\mathbf{A}_\omega^{\text{all}} = \begin{bmatrix} \mathbf{A}_\omega^{1 \rightarrow L-1} \\ \dots \\ \mathbf{A}_\omega^{1 \rightarrow 1} \end{bmatrix}, \quad \mathbf{b}_\omega^{\text{all}} = \begin{bmatrix} \mathbf{b}_\omega^{1 \rightarrow L-1} \\ \dots \\ \mathbf{b}_\omega^{1 \rightarrow 1} \end{bmatrix}, \quad \mathbf{A}_\omega^{\text{all}} \in \mathbb{R}^{(\prod_{\ell=1}^{L-1} D^\ell) \times S}, \quad \mathbf{b}_\omega^{\text{all}} \in \mathbb{R}^{\prod_{\ell=1}^{L-1} D^\ell}. \quad (12)$$

Corollary 2. *The \mathcal{H} -representation of the polyhedral region ω is given by*

$$\omega = \{\mathbf{z} \in \mathbb{R}^S : \mathbf{A}_\omega^{\text{all}} \mathbf{z} < -\mathbf{q}^{\text{all}}(\omega) \odot \mathbf{b}_\omega^{\text{all}}\} = \bigcap_{\ell=1}^{L-1} \{\mathbf{z} \in \mathbb{R}^S : \mathbf{A}_\omega^{1 \rightarrow \ell} \mathbf{z} < -\mathbf{q}^\ell(\omega) \odot \mathbf{b}_\omega^{1 \rightarrow \ell}\}, \quad (13)$$

with \odot the Hadamard product.

From the above, it is clear that the sign locates in which side of each hyperplane the region is located. We now have a direct way to obtain the polytope ω from its sign pattern $\mathbf{q}^{\text{all}}(\omega)$ or equivalently from an input $\mathbf{z} \in \omega$; the only task left is to obtain the entire partition Ω collecting all the DN regions, which we now propose to do via a simple scheme.

Partition Cells Enumeration. The search for all cells in a partition is known as the cell enumeration problem and has been extensively studied in the context of specific partitions such as hyperplane arrangements [54–56]. In our case however, the set of inequalities of different regions changes. In fact, for any neighbour region, not only the sign pattern \mathbf{q}^{all} will change but also $\mathbf{A}_\omega^{\text{all}}$ and $\mathbf{b}_\omega^{\text{all}}$ due to the

²looser condition can be put as the up-to-layer ℓ mapping is a CPA on a coarser partition than Ω but this is sufficient for our goal.

composition of layers. In fact, changing one activation state say -1 to 1 for a specific unit at layer ℓ will alter the affine parameters from (10) and (11) due to the layer composition. As such, we propose to enumerate all the cells $\omega \in \Omega$ with a deterministic algorithm that starts from an initial region and recursively explores its neighbouring cells until all have been visited while recomputing the inequality system at each step. To do so, consider the initial region ω_0 . First, one finds all the non-redundant inequalities of the inequality system (13), the remaining inequalities define the faces of the polytope ω . Second, one obtains any of the neighbouring regions sharing a face with ω_0 by switching the sign in the entry of $\mathbf{q}(\omega_0)$ corresponding to the considered face. Repeat this for all non-redundant inequalities to obtain all the adjacent regions to ω_0 sharing a face with it. Each altered code defines an adjacent region and its system of inequality can be obtained as per Lemma 2. Doing so for all the faces of the initial region and then iterating this process on all the newly discovered regions will enumerate the entire partition Ω . We summarize this in Algo 1 in the appendix and illustrate this recursive procedure in Fig. 1.

We now have each cell as a polytope and enumerated the partition Ω , we can now turn into the computation of the marginal and posterior DGN distributions.

B Analytical Moments for truncated Gaussian

To lighten the derivation, we introduce extend the $[\cdot]$ indexing operator such that for example for a matrix, $[\cdot]_{-k,\cdot}$ means that all the rows but the k^{th} are taken, and all columns are taken. Also, $[\cdot]_{(k,l),\cdot}$ means that only the k^{th} and l^{th} rows are taken and all the columns. Let also introduce the following quantities

$$\begin{aligned} [F(\mathbf{a}, \Sigma)]_k &= \phi([\mathbf{a}]_k; 0, [\Sigma]_{k,k}) \Phi_{[[\mathbf{a}]_{-k,\infty})}(\boldsymbol{\mu}(k), \Sigma(k)) \\ [G(\mathbf{a}, \Sigma)]_{k,l} &= \phi([\mathbf{a}]_{(k,l)}; 0, [\Sigma]_{(k,j),(k,j)}) \Phi_{[[\mathbf{a}]_{-(k,l),\infty})}(\boldsymbol{\mu}((k,l)), \Sigma((k,l))) \\ H(\mathbf{a}, \Sigma) &= G(\mathbf{a}, \Sigma) + \text{diag} \left(\frac{\mathbf{l} \odot F(\mathbf{l}, \Sigma) - (\Sigma \odot G(\mathbf{l}, \Sigma)) \mathbf{1}}{\text{diag}(\Sigma)} \right) \end{aligned}$$

with $\boldsymbol{\mu}(u) = [\Sigma]_{-u,u} [\Sigma]_{u,u}^{-1} [\mathbf{a}]_u$, and $\Sigma(u) = [\Sigma]_{-u,-u} - [\Sigma]_{-u,u} [\Sigma]_{u,u}^{-1} [\Sigma]_{-u,u}^T$. Thanks to the above form, we can now obtain the integral $e_\omega^0(\Sigma) \triangleq \Phi_\omega(\mathbf{0}, \Sigma)$ and the first two moments of a centered truncated gaussian $e_\omega^1(\Sigma) \triangleq \int_\omega \mathbf{z} \phi(\mathbf{z}; \mathbf{0}, \Sigma)$ and $E_\omega^2(\Sigma) \triangleq \int_\omega \mathbf{z} \mathbf{z}^T \phi(\mathbf{z}; \mathbf{0}, \Sigma)$

Corollary 3. *The integral and first two moments of a centered truncated gaussian are given by*

$$e_\omega^0(\Sigma) = \sum_{\Delta \in T(\omega)} \sum_{(s,C) \in T(\Delta)} s \Phi_{[\mathbf{l}(C), \infty)}(0, R_C \Sigma R_C^T) dz, \quad (14)$$

$$e_\omega^1(\Sigma) = \Sigma \sum_{\Delta \in T(\omega)} \sum_{(s,C) \in T(\Delta)} s R_C^T F(\mathbf{l}_{\omega,C}, R_C \Sigma R_C^T), \quad (15)$$

$$E_\omega^2(\Sigma) = \Sigma \left(\sum_{\Delta \in T(\omega)} \sum_{(s,C) \in T(\Delta)} s R_C^T (H(\mathbf{l}_{\omega,C}, R_C \Sigma R_C^T)) R_C \right) \Sigma + e_\omega^0(\Sigma) \Sigma \quad (16)$$

To simplify notations let consider the following notation of the posterior (6) where are incorporate the terms independent of \mathbf{z} into

$$\alpha_\omega(\mathbf{x}) = \frac{\phi(\mathbf{x}; B_\omega, \Sigma_\mathbf{x} + A_\omega \Sigma_z A_\omega^T)}{\sum_\omega \phi(\mathbf{x}; B_\omega, \Sigma_\mathbf{x} + A_\omega \Sigma_z A_\omega^T) \Phi_\omega(\boldsymbol{\mu}_\omega(\mathbf{x}), \Sigma_\omega)}, \quad (17)$$

leading to $p(\mathbf{z}|\mathbf{x}) = \sum_{\omega \in \Omega} \delta_\omega(\mathbf{z}) \alpha_\omega(\mathbf{x}) \phi(\mathbf{z}; \boldsymbol{\mu}_\omega(\mathbf{x}), \Sigma_\omega)$.

Theorem 2. *The first (per region) moments of the DGN posterior are given by*

$$\begin{aligned} \mathbb{E}_{\mathbf{z}|\mathbf{x}}[\mathbb{1}_{\mathbf{z} \in \omega}] &= \alpha_\omega(\mathbf{x}) e_\omega^0(\Sigma_\omega), \\ \mathbb{E}_{\mathbf{z}|\mathbf{x}}[\mathbf{z} \mathbb{1}_{\mathbf{z} \in \omega}] &= \alpha_\omega(\mathbf{x}) (e_{\omega - \boldsymbol{\mu}_\omega(\mathbf{x})}^1(\Sigma_\omega) + e_{\omega - \boldsymbol{\mu}_\omega(\mathbf{x})}^0(\Sigma_\omega) \boldsymbol{\mu}_\omega(\mathbf{x})) \end{aligned}$$

$$\begin{aligned} \mathbb{E}_{\mathbf{z}|\mathbf{x}}[\mathbf{z}\mathbf{z}^T \mathbb{1}_{\mathbf{z}\in\omega}] &= \alpha_\omega(\mathbf{x}) (\mathbf{E}_{\omega-\mu_\omega(\mathbf{x})}^2(\boldsymbol{\Sigma}_\omega) + e_{\omega-\mu_\omega(\mathbf{x})}^1(\boldsymbol{\Sigma}_\omega) \boldsymbol{\mu}_\omega(\mathbf{x})^T \\ &\quad + \boldsymbol{\mu}_{\omega-\mu_\omega(\mathbf{x})}(\mathbf{x}) e_{\omega-\mu_\omega(\mathbf{x})}^1(\mathbf{x})^T + \boldsymbol{\mu}_\omega(\mathbf{x}) \boldsymbol{\mu}_\omega(\mathbf{x})^T e_\omega^0(\mathbf{x})) \end{aligned}$$

which we denote $\mathbb{E}_{\mathbf{z}|\mathbf{x}}[\mathbb{1}_{\mathbf{z}\in\omega}] \triangleq e_\omega^0(\mathbf{x})$, $\mathbb{E}_{\mathbf{z}|\mathbf{x}}[\mathbf{z} \mathbb{1}_{\mathbf{z}\in\omega}] \triangleq e_\omega^1(\mathbf{x})$ and $\mathbb{E}_{\mathbf{z}|\mathbf{x}}[\mathbf{z}\mathbf{z}^T \mathbb{1}_{\mathbf{z}\in\omega}] \triangleq \mathbf{E}_\omega^2(\mathbf{x})$. (Proof in F.9.)

C Implementation Details

The Delaunay triangulation needs the \mathcal{V} -representation of ω , the vertices which convex hull form the region [57]. Given that we have the \mathcal{H} -representation, finding the vertices is known as the vertex enumeration problem [58]. To compute the triangulation we use the Python scipy [59] implementation which interfaces the C/C++ Qhull implementation [60]. To compute the $\mathcal{H} \mapsto \mathcal{V}$ representation and vice-versa we leverage pycddlib³ which interfaces the C/C++ cddlib library⁴ employing the double description method [61].

D Figures

We demonstrate here additional figures for the posterior and marginal distribution of a DGN.

³<https://pypi.org/project/pycddlib/>

⁴https://inf.ethz.ch/personal/fukudak/cdd_home/index.html

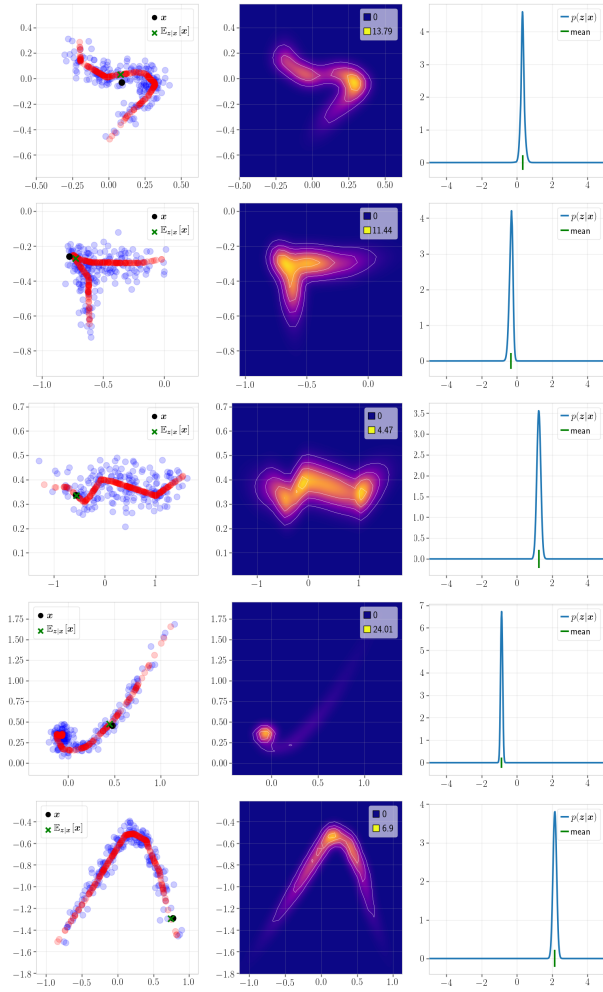


Figure 8: Additional random DGNs with their samples, the posterior and the marginal distributions.

E Algorithms

Algorithm 1: SearchRegion

Data: Starting region ω and $\mathbf{q}(\omega)$, initial set (Ω)

Result: Updated Ω

if $\omega \notin \Omega$ **then**

 | $\Omega \leftarrow \Omega \cup \{\omega\};$

else

 | Quit

end

$I = \text{reduce}(A_\omega^{\text{all}}, B_\omega^{\text{all}});$

for $i \in I$ **do**

 | SearchRegion(flip($\mathbf{q}(\omega)$), i), Ω);

end

F Proofs

In this section we provide all the proofs for the main paper theoretical claims. In particular we will go through the derivations of the per region posterior first moments and then the derivation of the expectation and maximization steps.

F.1 Proof of Lemma 1

Proof. The proof consists of expressing the conditional distribution and using the properties of DGN with piecewise affine nonlinearities. We are able to split the distribution into a mixture model as follows:

$$\begin{aligned}
 p(\mathbf{x}|\mathbf{z}) &= \frac{1}{(2\pi)^{D/2} \sqrt{|\det \Sigma_{\mathbf{x}}|}} e^{-\frac{1}{2}(\mathbf{x}-g(\mathbf{z}))^T \Sigma_{\mathbf{x}}^{-1}(\mathbf{x}-g(\mathbf{z}))} \\
 &= \frac{1}{(2\pi)^{D/2} \sqrt{|\det \Sigma_{\mathbf{x}}|}} e^{-\frac{1}{2}(\mathbf{x}-\sum_{\omega \in \Omega} \mathbb{1}_{\mathbf{z} \in \omega} (A_\omega \mathbf{z} + B_\omega))^T \Sigma_{\mathbf{x}}^{-1}(\mathbf{x}-\sum_{\omega \in \Omega} \mathbb{1}_{\mathbf{z} \in \omega} (A_\omega \mathbf{z} + B_\omega))} \\
 &= \frac{1}{(2\pi)^{D/2} \sqrt{|\det \Sigma_{\mathbf{x}}|}} e^{-\frac{1}{2} \sum_{\omega \in \Omega} \mathbb{1}_{\mathbf{z} \in \omega} (\mathbf{x} - (A_\omega \mathbf{z} + B_\omega))^T \Sigma_{\mathbf{x}}^{-1}(\mathbf{x} - (A_\omega \mathbf{z} + B_\omega))} \\
 &= \sum_{\omega \in \Omega} \mathbb{1}_{\mathbf{z} \in \omega} \frac{1}{(2\pi)^{D/2} \sqrt{|\det \Sigma_{\mathbf{x}}|}} e^{-\frac{1}{2}(\mathbf{x} - (A_\omega \mathbf{z} + B_\omega))^T \Sigma_{\mathbf{x}}^{-1}(\mathbf{x} - (A_\omega \mathbf{z} + B_\omega))} \\
 &= \sum_{\omega \in \Omega} \mathbb{1}_{\mathbf{z} \in \omega} \phi(\mathbf{x}|A_\omega \mathbf{z} + B_\omega, \Sigma_{\mathbf{x}})
 \end{aligned}$$

□

F.2 Proof of Proposition 1

Proof. This result is direct by noticing that the probability to obtain a specific region slope and bias is the probability that the sampled latent vector lies in the corresponding region. This probability is obtained simply by integrating the latent gaussian distribution on the region. We obtain the result of the proposition. □

F.3 Proof of Theorem 1

Proof. For the first part, we simply leverage the known result from linear Gaussian models [62] stating that

$$p(\mathbf{z}|\mathbf{x}) = \frac{p(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{p(\mathbf{x})}$$

$$\begin{aligned}
&= \frac{1}{p(\mathbf{x})} \frac{e^{-\frac{1}{2}(\mathbf{x}-g(\mathbf{z}))^T \Sigma_{\mathbf{x}}^{-1}(\mathbf{x}-g(\mathbf{z}))}}{(2\pi)^{D/2} \sqrt{|\det(\Sigma_{\mathbf{x}})|}} \frac{e^{-\frac{1}{2}(\mathbf{z}-\boldsymbol{\mu})^T \Sigma_{\mathbf{z}}^{-1}(\mathbf{z}-\boldsymbol{\mu})}}{(2\pi)^{S/2} \sqrt{|\det(\Sigma_{\mathbf{z}})|}} \\
&= \frac{1}{p(\mathbf{x})} \left(\sum_{\omega \in \Omega} \mathbb{1}_{\mathbf{z} \in \omega} \frac{e^{-\frac{1}{2}(\mathbf{x}-A_{\omega}\mathbf{z}-B_{\omega})^T \Sigma_{\mathbf{x}}^{-1}(\mathbf{x}-A_{\omega}\mathbf{z}-B_{\omega})}}{(2\pi)^{D/2} \sqrt{|\det(\Sigma_{\mathbf{x}})|}} \right) \frac{e^{-\frac{1}{2}(\mathbf{z}-\boldsymbol{\mu})^T \Sigma_{\mathbf{z}}^{-1}(\mathbf{z}-\boldsymbol{\mu})}}{(2\pi)^{S/2} \sqrt{|\det(\Sigma_{\mathbf{z}})|}} \\
&= \frac{1}{p(\mathbf{x})} \sum_{\omega \in \Omega} \mathbb{1}_{\mathbf{z} \in \omega} \frac{e^{-\frac{1}{2}(\mathbf{x}-A_{\omega}\mathbf{z}-B_{\omega})^T \Sigma_{\mathbf{x}}^{-1}(\mathbf{x}-A_{\omega}\mathbf{z}-B_{\omega}) - \frac{1}{2}\mathbf{z}^T \Sigma_{\mathbf{z}}^{-1}\mathbf{z}}}{(2\pi)^{(S+D)/2} \sqrt{|\det(\Sigma_{\mathbf{x}})| |\det(\Sigma_{\mathbf{z}})|}} \\
&= \frac{1}{p(\mathbf{x})} \sum_{\omega \in \Omega} \mathbb{1}_{\mathbf{z} \in \omega} \frac{e^{-\frac{1}{2}((\mathbf{x}-B_{\omega})-A_{\omega}\mathbf{z})^T \Sigma_{\mathbf{x}}^{-1}((\mathbf{x}-B_{\omega})-A_{\omega}\mathbf{z}) - \frac{1}{2}\mathbf{z}^T \Sigma_{\mathbf{z}}^{-1}\mathbf{z}}}{(2\pi)^{(S+D)/2} \sqrt{|\det(\Sigma_{\mathbf{x}})| |\det(\Sigma_{\mathbf{z}})|}} \\
&= \frac{1}{p(\mathbf{x})} \sum_{\omega \in \Omega} \mathbb{1}_{\mathbf{z} \in \omega} \frac{e^{-\frac{1}{2}((A_{\omega}^T \Sigma_{\mathbf{x}}^{-1} A_{\omega} + \Sigma_{\mathbf{x}}^{-1})^{-1} A_{\omega}^T \Sigma_{\mathbf{x}}^{-1} (\mathbf{x}-B_{\omega}) - \mathbf{z})^T (A_{\omega}^T \Sigma_{\mathbf{x}}^{-1} A_{\omega} + \Sigma_{\mathbf{x}}^{-1}) ((A_{\omega}^T \Sigma_{\mathbf{x}}^{-1} A_{\omega} + \Sigma_{\mathbf{x}}^{-1})^{-1} A_{\omega}^T \Sigma_{\mathbf{x}}^{-1} (\mathbf{x}-B_{\omega}) - \mathbf{z})}}{(2\pi)^{(S+D)/2} \sqrt{|\det(\Sigma_{\mathbf{x}})| |\det(\Sigma_{\mathbf{z}})|}} \\
&\quad \times e^{-\frac{1}{2}((\mathbf{x}-B_{\omega})^T \Sigma_{\mathbf{x}}^{-1} (\mathbf{x}-B_{\omega})) + \frac{1}{2}((\mathbf{x}-B_{\omega})^T \Sigma_{\mathbf{x}}^{-1} A_{\omega} (A_{\omega}^T \Sigma_{\mathbf{x}}^{-1} A_{\omega} + \Sigma_{\mathbf{x}}^{-1})^{-1} A_{\omega}^T \Sigma_{\mathbf{x}}^{-1} (\mathbf{x}-B_{\omega}))} \\
&= \frac{1}{p(\mathbf{x})} \sum_{\omega \in \Omega} \mathbb{1}_{\mathbf{z} \in \omega} \frac{e^{-\frac{1}{2}((A_{\omega}^T \Sigma_{\mathbf{x}}^{-1} A_{\omega} + \Sigma_{\mathbf{x}}^{-1})^{-1} A_{\omega}^T \Sigma_{\mathbf{x}}^{-1} (\mathbf{x}-B_{\omega}) - \mathbf{z})^T (A_{\omega}^T \Sigma_{\mathbf{x}}^{-1} A_{\omega} + \Sigma_{\mathbf{x}}^{-1}) ((A_{\omega}^T \Sigma_{\mathbf{x}}^{-1} A_{\omega} + \Sigma_{\mathbf{x}}^{-1})^{-1} A_{\omega}^T \Sigma_{\mathbf{x}}^{-1} (\mathbf{x}-B_{\omega}) - \mathbf{z})}}{(2\pi)^{(S+D)/2} \sqrt{|\det(\Sigma_{\mathbf{x}})| |\det(\Sigma_{\mathbf{z}})|}} \\
&\quad \times e^{-\frac{1}{2}((\mathbf{x}-B_{\omega})^T (\Sigma_{\mathbf{x}}^{-1} - \Sigma_{\mathbf{x}}^{-1} A_{\omega} (A_{\omega}^T \Sigma_{\mathbf{x}}^{-1} A_{\omega} + \Sigma_{\mathbf{x}}^{-1})^{-1} A_{\omega}^T \Sigma_{\mathbf{x}}^{-1}) (\mathbf{x}-B_{\omega}))} \\
&= \frac{1}{p(\mathbf{x})} \sum_{\omega \in \Omega} \mathbb{1}_{\mathbf{z} \in \omega} \frac{e^{-\frac{1}{2}((A_{\omega}^T \Sigma_{\mathbf{x}}^{-1} A_{\omega} + \Sigma_{\mathbf{x}}^{-1})^{-1} A_{\omega}^T \Sigma_{\mathbf{x}}^{-1} (\mathbf{x}-B_{\omega}) - \mathbf{z})^T (A_{\omega}^T \Sigma_{\mathbf{x}}^{-1} A_{\omega} + \Sigma_{\mathbf{x}}^{-1}) ((A_{\omega}^T \Sigma_{\mathbf{x}}^{-1} A_{\omega} + \Sigma_{\mathbf{x}}^{-1})^{-1} A_{\omega}^T \Sigma_{\mathbf{x}}^{-1} (\mathbf{x}-B_{\omega}) - \mathbf{z})}}{(2\pi)^{(S+D)/2} \sqrt{|\det(\Sigma_{\mathbf{x}})| |\det(\Sigma_{\mathbf{z}})|}} \\
&\quad \times e^{-\frac{1}{2}((\mathbf{x}-B_{\omega})^T (\Sigma_{\mathbf{x}} + A_{\omega} \Sigma_{\mathbf{z}} A_{\omega}^T)^{-1} (\mathbf{x}-B_{\omega}))} \\
&= \frac{1}{p(\mathbf{x})} \sum_{\omega \in \Omega} \mathbb{1}_{\mathbf{z} \in \omega} \frac{e^{-\frac{1}{2}(\boldsymbol{\mu}_{\omega}(\mathbf{x}) - \mathbf{z})^T \Sigma_{\omega}^{-1}(\boldsymbol{\mu}_{\omega}(\mathbf{x}) - \mathbf{z})}}{(2\pi)^{(S+D)/2} \sqrt{|\det(\Sigma_{\mathbf{x}})| |\det(\Sigma_{\mathbf{z}})|}} e^{-\frac{1}{2}((\mathbf{x}-B_{\omega})^T (\Sigma_{\mathbf{x}} + A_{\omega} \Sigma_{\mathbf{z}} A_{\omega}^T)^{-1} (\mathbf{x}-B_{\omega}))}
\end{aligned}$$

with $\boldsymbol{\mu}_{\omega}(\mathbf{x}) = \Sigma_{\omega} A_{\omega}^T \Sigma_{\mathbf{x}}^{-1} (\mathbf{x} - B_{\omega})$ and $\Sigma_{\omega} = (A_{\omega}^T \Sigma_{\mathbf{x}}^{-1} A_{\omega} + \Sigma_{\mathbf{x}}^{-1})^{-1}$ as a result it corresponds to a mixture of truncated gaussian, each living on ω . Now we determine the renormalization constant:

$$\begin{aligned}
p(\mathbf{x}) &= \int p(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z} \\
&= \sum_{\omega \in \Omega} \int_{\omega} \mathbb{1}_{\mathbf{z} \in \omega} \frac{e^{-\frac{1}{2}(\boldsymbol{\mu}_{\omega}(\mathbf{x}) - \mathbf{z})^T \Sigma_{\omega}^{-1}(\boldsymbol{\mu}_{\omega}(\mathbf{x}) - \mathbf{z})}}{(2\pi)^{(S+D)/2} \sqrt{|\det(\Sigma_{\mathbf{x}})| |\det(\Sigma_{\mathbf{z}})|}} e^{-\frac{1}{2}((\mathbf{x}-B_{\omega})^T (\Sigma_{\mathbf{x}} + A_{\omega} \Sigma_{\mathbf{z}} A_{\omega}^T)^{-1} (\mathbf{x}-B_{\omega}))} d\mathbf{z} \\
&= \sum_{\omega \in \Omega} \mathbb{1}_{\mathbf{z} \in \omega} \frac{e^{-\frac{1}{2}((\mathbf{x}-B_{\omega})^T (\Sigma_{\mathbf{x}} + A_{\omega} \Sigma_{\mathbf{z}} A_{\omega}^T)^{-1} (\mathbf{x}-B_{\omega}))}}{(2\pi)^{D/2} \sqrt{|\det(\Sigma_{\mathbf{x}})| |\det(\Sigma_{\mathbf{z}})|}} \sqrt{\det(\Sigma_{\omega})} \int_{\omega} \phi(\mathbf{z}; \boldsymbol{\mu}_{\omega}(\mathbf{x}), \Sigma_{\omega}) d\mathbf{z} \\
&= \sum_{\omega \in \Omega} \mathbb{1}_{\mathbf{z} \in \omega} \frac{e^{-\frac{1}{2}((\mathbf{x}-B_{\omega})^T (\Sigma_{\mathbf{x}} + A_{\omega} \Sigma_{\mathbf{z}} A_{\omega}^T)^{-1} (\mathbf{x}-B_{\omega}))}}{(2\pi)^{D/2} \sqrt{|\det(\Sigma_{\mathbf{x}})| |\det(\Sigma_{\mathbf{z}})|}} \sqrt{\det(\Sigma_{\omega})} \Phi_{\omega}(\boldsymbol{\mu}_{\omega}(\mathbf{x}), \Sigma_{\omega}) \\
&= \sum_{\omega \in \Omega} \mathbb{1}_{\mathbf{z} \in \omega} \frac{\sqrt{\det(\Sigma_{\mathbf{x}} + A_{\omega} \Sigma_{\mathbf{z}} A_{\omega}^T) \det(\Sigma_{\omega})}}{\sqrt{|\det(\Sigma_{\mathbf{x}})| |\det(\Sigma_{\mathbf{z}})|}} \phi(\mathbf{x}; B_{\omega}, \Sigma_{\mathbf{x}} + A_{\omega} \Sigma_{\mathbf{z}} A_{\omega}^T) \Phi_{\omega}(\boldsymbol{\mu}_{\omega}(\mathbf{x}), \Sigma_{\omega}),
\end{aligned}$$

now using the Matrix determinant lemma [63] we have that $\det(\Sigma_{\mathbf{x}} + A_{\omega} \Sigma_{\mathbf{z}} A_{\omega}^T) = \det(\Sigma_{\mathbf{z}}^{-1} + A_{\omega}^T \Sigma_{\mathbf{x}}^{-1} A_{\omega}) \det(\Sigma_{\mathbf{x}}) \det(\Sigma_{\mathbf{z}})$ leading to

$$p(\mathbf{x}) = \sum_{\omega} \phi(\mathbf{x}; B_{\omega}, \Sigma_{\mathbf{x}} + A_{\omega} \Sigma_{\mathbf{z}} A_{\omega}^T) \Phi_{\omega}(\boldsymbol{\mu}_{\omega}(\mathbf{x}), \Sigma_{\omega}),$$

$$p(\mathbf{z}|\mathbf{x}) = \sum_{\omega} \delta_{\omega}(\mathbf{z}) \frac{\phi(\mathbf{x}; B_{\omega}, \Sigma_{\mathbf{x}} + A_{\omega} \Sigma_{\mathbf{z}} A_{\omega}^T) \phi(\mathbf{z}; \boldsymbol{\mu}_{\omega}(\mathbf{x}), \Sigma_{\omega})}{\sum_{\omega} \phi(\mathbf{x}; B_{\omega}, \Sigma_{\mathbf{x}} + A_{\omega} \Sigma_{\mathbf{z}} A_{\omega}^T) \Phi_{\omega}(\boldsymbol{\mu}_{\omega}(\mathbf{x}), \Sigma_{\omega})}.$$

□

F.4 Proof of Lemma 2

The proof will consist of observing that the posterior (prior rewriting) can be expressed as a softmax of a quantity rescaled by the standard deviation.

Proof.

$$\begin{aligned} \log(\phi(\mathbf{z}; \boldsymbol{\mu}_{\omega}(\mathbf{x}), \Sigma_{\omega})) &= -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_{\omega})^T \Sigma_{\omega}^{-1}(\mathbf{x} - \boldsymbol{\mu}_{\omega}) - \frac{1}{2} \log(\det(\Sigma_{\omega})) + cst \\ &= -\frac{1}{2}(\mathbf{z} - \Sigma_{\omega} A_{\omega}^T \Sigma_{\mathbf{x}}^{-1} (A_0 \mathbf{z}_0 + B_0 - B_{\omega}))^T \Sigma_{\omega}^{-1} (\mathbf{z} - \Sigma_{\omega} A_{\omega}^T \Sigma_{\mathbf{x}}^{-1} (A_0 \mathbf{z}_0 + B_0 - B_{\omega}))^T \\ &\quad - \frac{1}{2} \log(\det(\Sigma_{\omega})) + cst \\ &= -\frac{1}{2}(\mathbf{z} - (A_{\omega}^T \Sigma_{\mathbf{x}}^{-1} A_{\omega})^{-1} A_{\omega}^T \Sigma_{\mathbf{x}}^{-1} (A_0 \mathbf{z}_0 + B_0 - B_{\omega}))^T \Sigma_{\omega}^{-1} (\mathbf{z} - (A_{\omega}^T \Sigma_{\mathbf{x}}^{-1} A_{\omega})^{-1} A_{\omega}^T \Sigma_{\mathbf{x}}^{-1} (A_0 \mathbf{z}_0 + B_0 - B_{\omega}))^T \\ &\quad - \frac{1}{2} \log(\det(\Sigma_{\omega})) + cst \end{aligned}$$

where we used the following result to develop $\boldsymbol{\mu}_{\omega}(\mathbf{x})$

$$\begin{aligned} \Sigma_{\omega} &= (A_{\omega}^T \Sigma_{\mathbf{x}}^{-1} A_{\omega} + \Sigma_{\mathbf{z}}^{-1})^{-1} = (A_{\omega}^T \Sigma_{\mathbf{x}}^{-1} A_{\omega} + (A_{\omega}^T \Sigma_{\mathbf{x}}^{-1} A_{\omega})(A_{\omega}^T \Sigma_{\mathbf{x}}^{-1} A_{\omega})^{-1} \Sigma_{\mathbf{z}}^{-1})^{-1} \\ &= (A_{\omega}^T \Sigma_{\mathbf{x}}^{-1} A_{\omega})^{-1} (I + (A_{\omega}^T \Sigma_{\mathbf{x}}^{-1} A_{\omega})^{-1} \Sigma_{\mathbf{z}}^{-1})^{-1} \\ &= (A_{\omega}^T \Sigma_{\mathbf{x}}^{-1} A_{\omega})^{-1} \text{ as } (\Sigma_{\mathbf{z}} A_{\omega}^T \Sigma_{\mathbf{x}}^{-1} A_{\omega})^{-1} \rightarrow \mathbf{0}. \end{aligned}$$

if we are in the same region ω than \mathbf{z}_0 then the above becomes

$$\arg \max_{\mathbf{z} \in \omega_0} \log(\phi(\mathbf{z}; \boldsymbol{\mu}_{\omega}(\mathbf{x}), \Sigma_{\omega})) = \arg \max_{\mathbf{z} \in \omega_0} -\frac{1}{2}(\mathbf{z} - \mathbf{z}_0)^T \Sigma_{\omega}^{-1} (\mathbf{z} - \mathbf{z}_0) = \mathbf{z}_0,$$

and since we know that we are in the same region, the argmax $\mathbf{z} = \mathbf{z}_0$ lies in this region and thus is the maximum of the posterior.

□

F.5 Proof of Lemma 3

Proof. The sign vectors represent the sign of each pre-activation feature maps. The key here is that when changing the sign, the input passes through the knot of the corresponding activation function of that layer. This implies a change in the region in the DGN input space. In fact, without degenerate weights and with nonzero activation functions, a change in any dimension of the sign vector (used to form the per region slope and bias) impact a change in the affine mapping used to map inputs \mathbf{z} to outputs \mathbf{x} . As such, whenever a sign changes, the affine mapping changes, leading to a change of region in the DGN input space. As the sign vector is formed from the DGN input space, and we restrict ourselves to the image of this mapping, there does not exist a sign pattern/configuration not reachable by the DGN (otherwise it would not be in the image of this mapping). Now for the other inclusion, recall that a change in region and thus in per region affine mapping can only occur with a change of pre-activation sign pattern. □

F.6 Proof of Corollary 2

Proof. From the above result, it is clear that the preactivation roots define the boundaries of the regions. Obtaining the hyperplane representation of the region thus simply consists of reexpressing this statement with the explicit pre-activation hyperplanes for all the layers and units, the intersection between layers coming from the subdivision. For additional details please see [25]. □

F.7 Proof of Lemma 4

Proof. The proof consists of rearranging the terms from the inclusion-exclusion formula as in

$$\begin{aligned}
& \sum_{J \subseteq \{1, \dots, F\}, J \neq \emptyset} (-1)^{|J|+1} (\cap_{j \in J} A_j) = \cup_i A_i \\
(-1)^{F+1} S + & \sum_{J \subseteq \{1, \dots, F\}, J \neq \emptyset, |J| < F} (-1)^{|J|+1} (\cap_{j \in J} A_j) = \cup_i A_i \\
(-1)^{F+1} S = & \cup_i A_i - \sum_{J \subseteq \{1, \dots, F\}, J \neq \emptyset, |J| < F} (-1)^{|J|+1} (\cap_{j \in J} A_j) \\
S = & (-1)^{F+1} \cup_i A_i - (-1)^{F+1} \sum_{J \subseteq \{1, \dots, F\}, J \neq \emptyset, |J| < F} (-1)^{|J|+1} (\cap_{j \in J} A_j) \\
S = & (-1)^{F+1} \cup_i A_i + \sum_{J \subseteq \{1, \dots, F\}, J \neq \emptyset, |J| < F} (-1)^{|J|+1+F} (\cap_{j \in J} A_j)
\end{aligned}$$

then by application of Chasles rule [64], the integral domain can be decomposed into the signed sum of per cone integration. Finally, a simplex in dimension S has $S + 1$ faces, making $F = S + 1$ and leading to the desired result. \square

F.8 Proof of Moments

Lemma 6. *The first moments of Gaussian integration on an open rectangle defined by its lower limits \mathbf{a} is given by*

$$\int_{\mathbf{a}}^{\infty} \mathbf{z} \phi(\mathbf{0}, \Sigma) d\mathbf{z} = \Sigma F(\mathbf{a}), \tag{18}$$

$$\int_{\mathbf{a}}^{\infty} \mathbf{z} \mathbf{z}^T \phi(\mathbf{0}, \Sigma) d\mathbf{z} = \Phi_{[\mathbf{a}, \infty)}(\mathbf{0}, \Sigma) \Sigma + \Sigma \left(G(\mathbf{a}) + \frac{\mathbf{a} \odot F(\mathbf{a}) - (\Sigma \odot G(\mathbf{a})) \mathbf{1}}{\text{diag}(\Sigma)} \right) \Sigma. \tag{19}$$

where the division is performed elementwise.

Proof. First moment:

$$\begin{aligned}
\int_{\omega} \mathbf{z} \phi(\mathbf{x}; \mathbf{0}, \Sigma) d\mathbf{z} &= \int_{\omega} \mathbf{z} \frac{e^{-\frac{1}{2} \mathbf{z}^T \Sigma^{-1} \mathbf{z}}}{(2\pi)^{K/2} |\det(\Sigma)|^{1/2}} d\mathbf{z} \\
&= \sum_{\Delta \in S(\omega)} \sum_{(s, C) \in T(\Delta)} s \int_C \mathbf{z} \frac{e^{-\frac{1}{2} (R_C \mathbf{z})^T (R_C^T)^{-1} \Sigma_{\omega}^{-1} R_C^{-1} R_C \mathbf{z}}}{(2\pi)^{K/2} |\det(\Sigma_{\omega})|^{1/2}} d\mathbf{z} \\
&= \sum_{\Delta \in S(\omega)} \sum_{(s, C) \in T(\Delta)} s \int_{\mathbf{u}(C)} R^{-1} \mathbf{u} \frac{e^{-\frac{1}{2} \mathbf{u}^T (R_C \Sigma_{\omega} R_C^T)^{-1} \mathbf{u}}}{(2\pi)^{K/2} |\det(R_C)| |\det(\Sigma_{\omega})|^{1/2}} d\mathbf{u} \\
&= \sum_{\Delta \in S(\omega)} \sum_{(s, C) \in T(\Delta)} s R_C^{-1} \int_{\mathbf{u}(C)} \mathbf{u} \phi(\mathbf{u}; \mathbf{0}, R_C \Sigma_{\omega} R_C^T) d\mathbf{u} \\
&= \sum_{\Delta \in S(\omega)} \sum_{(s, C) \in T(\Delta)} s R_C^{-1} (R_C \Sigma_{\omega} R_C^T F(\mathbf{u}(C))) \\
&= \Sigma_{\omega} \sum_{\Delta \in S(\omega)} \sum_{(s, C) \in T(\Delta)} s R_C^T F(\mathbf{u}(C))
\end{aligned}$$

Second moment

$$\begin{aligned}
\int_{\omega} \mathbf{z}\mathbf{z}^T \phi(\mathbf{x}; \mathbf{0}, \Sigma) d\mathbf{z} &= \int_{\omega} \mathbf{z}\mathbf{z}^T \frac{e^{-\frac{1}{2}\mathbf{z}^T \Sigma^{-1} \mathbf{z}}}{(2\pi)^{K/2} |\det(\Sigma)|^{1/2}} d\mathbf{z} \\
&= \sum_{\Delta \in S(\omega)} \sum_{(s,C) \in T(\Delta)} s \int_C \mathbf{z}\mathbf{z}^T \frac{e^{-\frac{1}{2}(R_C \mathbf{y})^T (R_C^T)^{-1} \Sigma_{\omega}^{-1} R_C^{-1} R_C \mathbf{y}}}{(2\pi)^{K/2} |\det(\Sigma_{\omega})|^{1/2}} d\mathbf{z} \\
&= \sum_{\Delta \in S(\omega)} \sum_{(s,C) \in T(\Delta)} s \int_{\mathbf{l}(C)} R_C^{-1} \mathbf{u}\mathbf{u}^T (R_C^{-1})^T \frac{e^{-\frac{1}{2}\mathbf{u}^T (R_C \Sigma_{\omega} R_C^T)^{-1} \mathbf{u}}}{(2\pi)^{K/2} |\det(R_C)| |\det(\Sigma_{\omega})|^{1/2}} d\mathbf{u} \\
&= \sum_{\Delta \in S(\omega)} \sum_{(s,C) \in T(\Delta)} s R_C^{-1} \int_{\mathbf{l}(C)} \mathbf{u}\mathbf{u}^T \phi(\mathbf{u}; \mathbf{0}, R_C \Sigma_{\omega} R_C^T) d\mathbf{u} (R_C^{-1})^T \\
&= \sum_{\Delta \in S(\omega - \boldsymbol{\mu}_{\omega}(\mathbf{x}))} \sum_{(s,C) \in T(\Delta)} s R_C^{-1} \left[\Phi_{[\mathbf{l}(C), \infty)}(\mathbf{0}, R_C \Sigma_{\omega} R_C^T) R_C \Sigma_{\omega} R_C^T \right. \\
&\quad \left. + R_C \Sigma_{\omega} R_C^T \left(\frac{\mathbf{l}(C) \odot F(\mathbf{l}(C)) + (R_C \Sigma_{\omega} R_C^T \odot G(\mathbf{l}(C))) \mathbf{1}}{\text{diag}(R_C \Sigma_{\omega} R_C^T)} \right) (R_C \Sigma_{\omega} R_C^T)^T \right] (R_C^{-1})^T \\
&= \sum_{\Delta \in S(\omega - \boldsymbol{\mu}_{\omega}(\mathbf{x}))} \sum_{(s,C) \in T(\Delta)} s \left[\Phi_{[\mathbf{l}(C), \infty)}(\mathbf{0}, R_C \Sigma_{\omega} R_C^T) \Sigma_{\omega} \right. \\
&\quad \left. + \Sigma_{\omega} R_C^T \left(\frac{\mathbf{l}(C) \odot F(\mathbf{l}(C)) + (R_C \Sigma_{\omega} R_C^T \odot G(\mathbf{l}(C))) \mathbf{1}}{\text{diag}(R_C \Sigma_{\omega} R_C^T)} \right) R_C \Sigma_{\omega} \right] \\
&= e_{\omega - \boldsymbol{\mu}_{\omega}(\mathbf{x})}^0 \Sigma_{\omega} + \Sigma_{\omega} \left[\sum_{\Delta \in S(\omega - \boldsymbol{\mu}_{\omega}(\mathbf{x}))} \sum_{(s,C) \in T(\Delta)} s R_C^T \left(\frac{\mathbf{l}(C) \odot F(\mathbf{l}(C)) + (R_C \Sigma_{\omega} R_C^T \odot G(\mathbf{l}(C))) \mathbf{1}}{\text{diag}(R_C \Sigma_{\omega} R_C^T)} \right) R_C \right] \Sigma_{\omega}
\end{aligned}$$

□

F.9 Proof of Theorem 2

Proof. Constant: $R_C = \begin{pmatrix} C^T \\ H^T \Sigma_{\omega}^{-1} \end{pmatrix}$

$$\int_{\omega} p(\mathbf{z}|\mathbf{x}) d\mathbf{z} = \alpha_{\omega}(\mathbf{x}) \int_{\omega} \phi(\mathbf{z}; \boldsymbol{\mu}_{\omega}(\mathbf{x}), \Sigma_{\omega}) d\mathbf{z} = \alpha_{\omega}(\mathbf{x}) \int_{\omega - \boldsymbol{\mu}_{\omega}(\mathbf{x})} \phi(\mathbf{z}; \mathbf{0}, \Sigma_{\omega}) d\mathbf{z} = \alpha_{\omega}(\mathbf{x}) e_{\omega - \boldsymbol{\mu}_{\omega}(\mathbf{x})}^0$$

First moment:

$$\begin{aligned}
\int_{\omega} \mathbf{z} p(\mathbf{z}|\mathbf{x}) d\mathbf{z} &= \alpha_{\omega}(\mathbf{x}) \int_{\omega} \mathbf{z} \frac{e^{-\frac{1}{2}(\mathbf{z} - \boldsymbol{\mu}_{\omega}(\mathbf{x}))^T \Sigma_{\omega}^{-1} (\mathbf{z} - \boldsymbol{\mu}_{\omega}(\mathbf{x}))}}{(2\pi)^{K/2} |\det(\Sigma_{\omega})|^{1/2}} d\mathbf{z} \\
&= \alpha_{\omega}(\mathbf{x}) \int_{\omega - \boldsymbol{\mu}_{\omega}(\mathbf{x})} (\mathbf{y} + \boldsymbol{\mu}_{\omega}(\mathbf{x})) \frac{e^{-\frac{1}{2}\mathbf{y}^T \Sigma_{\omega}^{-1} \mathbf{y}}}{(2\pi)^{K/2} |\det(\Sigma_{\omega})|^{1/2}} d\mathbf{y} \\
&= \alpha_{\omega}(\mathbf{x}) \left(e_{\omega - \boldsymbol{\mu}_{\omega}(\mathbf{x})}^1 + e_{\omega - \boldsymbol{\mu}_{\omega}(\mathbf{x})}^0 \boldsymbol{\mu}_{\omega}(\mathbf{x}) \right)
\end{aligned}$$

Second moment:

$$\int_{\omega} \mathbf{z}\mathbf{z}^T p(\mathbf{z}|\mathbf{x}) d\mathbf{z} = \alpha_{\omega}(\mathbf{x}) \int_{\omega} \mathbf{z}\mathbf{z}^T \phi(\mathbf{z}; \boldsymbol{\mu}_{\omega}(\mathbf{x}), \Sigma_{\omega}) d\mathbf{z}$$

$$\begin{aligned}
&= \alpha_\omega(\mathbf{x}) \int_{\omega - \mu_\omega(\mathbf{x})} (\mathbf{y} + \mu_\omega(\mathbf{x})) (\mathbf{y} + \mu_\omega(\mathbf{x}))^T \phi(\mathbf{z}; \mathbf{0}, \Sigma_\omega) d\mathbf{z} \\
&= \alpha_\omega(\mathbf{x}) \left(\mathbf{E}^2 + \mu_\omega(\mathbf{x}) e_\omega^1(\Sigma_\omega)^T + e_\omega^1(\Sigma_\omega) \mu_\omega(\mathbf{x})^T + \mu_\omega(\mathbf{x}) \mu_\omega(\mathbf{x})^T e_{\omega - \mu_\omega(\mathbf{x})}^0(\Sigma_\omega) \right)
\end{aligned}$$

□

G Proof of EM-step

We now derive the expectation maximization steps for a piecewise affine and continuous DGN.

G.1 E-step derivation

$$E_{\mathbf{z}|\mathbf{x}}[(A_\omega \mathbf{z} + B_\omega) \mathbb{1}_\omega] = A m_\omega^1 + B e_\omega^0 \quad (20)$$

$$E_{\mathbf{z}|\mathbf{x}}[\mathbf{z}^T A_\omega^T A_\omega \mathbf{z} \mathbb{1}_\omega] = \text{Tr}(A_\omega^T A_\omega m^2) \quad (21)$$

$$\begin{aligned}
E_{Z|X} [\log(p_{X|Z}(\mathbf{x}|\mathbf{z})p_Z(\mathbf{z}))] &= E_{Z|X} \left[\log \left(\frac{e^{-\frac{1}{2}(\mathbf{x}-g(\mathbf{z}))^T \Sigma_x^{-1}(\mathbf{x}-g(\mathbf{z}))}}{(2\pi)^{D/2} \sqrt{|\det(\Sigma_x)|}} \frac{e^{-\frac{1}{2}\mathbf{z}^T \Sigma_z^{-1} \mathbf{z}}}{(2\pi)^{S/2} \sqrt{|\det(\Sigma_z)|}} \right) \right] \\
&= -\log \left((2\pi)^{(S+D)/2} \sqrt{|\det(\Sigma_z)|} \sqrt{|\det(\Sigma_x)|} \right) - \frac{1}{2} E_{Z|X} \left[(\mathbf{x} - g(\mathbf{z}))^T \Sigma_x^{-1} (\mathbf{x} - g(\mathbf{z})) + \mathbf{z}^T \Sigma_z^{-1} \mathbf{z} \right] \\
&= -\log \left((2\pi)^{(S+D)/2} \sqrt{|\det(\Sigma_z)|} \sqrt{|\det(\Sigma_x)|} \right) \\
&\quad - \frac{1}{2} \left(\mathbf{x}^T \Sigma_x^{-1} \mathbf{x} + E_{Z|X} \left[-2\mathbf{x}^T \Sigma_x^{-1} g(\mathbf{z}) + g(\mathbf{z})^T \Sigma_x^{-1} g(\mathbf{z}) + \mathbf{z}^T \Sigma_z^{-1} \mathbf{z} \right] \right) \\
&= -\log \left((2\pi)^{(S+D)/2} \sqrt{|\det(\Sigma_z)|} \sqrt{|\det(\Sigma_x)|} \right) - \frac{1}{2} \left(\mathbf{x}^T \Sigma_x^{-1} \mathbf{x} + \text{Tr}(E_{Z|X}[\mathbf{z} \mathbf{z}^T \Sigma_z^{-1}]) \right. \\
&\quad \left. + E_{Z|X} \left[-2\mathbf{x}^T \Sigma_x^{-1} g(\mathbf{z}) + g(\mathbf{z})^T \Sigma_x^{-1} g(\mathbf{z}) \right] \right) \\
&= -\log \left((2\pi)^{(S+D)/2} \sqrt{|\det(\Sigma_z)|} \sqrt{|\det(\Sigma_x)|} \right) - \frac{1}{2} \left(\mathbf{x}^T \Sigma_x^{-1} \mathbf{x} - 2\mathbf{x}^T \Sigma_x^{-1} \left(\sum_\omega \mathbf{A}_\omega e_\omega^1(\mathbf{x}) + \mathbf{b}_\omega e_\omega^0(\mathbf{x}) \right) \right. \\
&\quad \left. + \sum_\omega e_\omega^0 \mathbf{b}_\omega^T \Sigma_x^{-1} \mathbf{b}_\omega + \text{Tr}(\mathbf{A}_\omega^T \Sigma_x^{-1} \mathbf{A}_\omega \mathbf{E}_\omega^2(\mathbf{x})) + 2(\mathbf{A}_\omega m_\omega^1(\mathbf{x}))^T \Sigma_x^{-1} \mathbf{b}_\omega + \text{Tr}(\Sigma_z^{-1} \mathbf{E}^2(\mathbf{x})) \right)
\end{aligned}$$

G.2 Proof of M step

Let first introduce some notations:

$$\mathbf{A}_\omega^{L \rightarrow i} \triangleq (\mathbf{A}_\omega^{L \rightarrow i})^T \text{ (back-propagation matrix to layer } i\text{),}$$

$$r_\omega^\ell(\mathbf{x}) \triangleq \left(\mathbf{x} e_\omega^0(\mathbf{x}) - \left(\mathbf{A}_\omega e_\omega^1(\mathbf{x}) + \sum_{i \neq \ell} m_\omega^0(\mathbf{x}) \mathbf{A}_\omega^{i+1 \rightarrow L} \mathbf{D}_\omega^i \mathbf{v}^i \right) \right) \text{ (expected residual without } \mathbf{v}^\ell\text{)}$$

$$\hat{\mathbf{z}}_\omega^\ell(\mathbf{x}) \triangleq \mathbf{D}_\omega^{\ell-1} (\mathbf{A}_\omega^{1 \rightarrow \ell-1} m_\omega^1(\mathbf{x}) + \mathbf{b}_\omega^{1 \rightarrow \ell-1} e_\omega^0) \text{ (expected feature map of layer } \ell\text{)}$$

we can now provide the analytical forms of the M step for each of the learnable parameters:

$$\Sigma_x^* = \frac{1}{N} \sum_{\mathbf{x}} \left(\mathbf{x} \mathbf{x}^T + \sum_\omega \mathbf{b}_\omega (\mathbf{b}_\omega m_\omega^0(\mathbf{x}) + 2\mathbf{A}_\omega e_\omega^1(\mathbf{x}))^T - 2\mathbf{x} (\hat{\mathbf{z}}_\omega^L(\mathbf{x}))^T + \mathbf{A}_\omega \mathbf{E}_\omega^2(\mathbf{x}) \mathbf{A}_\omega^T \right), \quad (22)$$

$$\mathbf{v}^{\ell*} = \left(\sum_{\mathbf{x}} \sum_{\omega} \mathbf{D}_{\omega}^{\ell} \mathbf{A}_{\omega}^{L \rightarrow \ell+1} \boldsymbol{\Sigma}_{\mathbf{x}}^{-1} \mathbf{A}_{\omega}^{\ell+1 \rightarrow L} \mathbf{D}_{\omega}^{\ell} \right)^{-1} \left(\sum_{\mathbf{x}} \sum_{\omega \in \Omega} \underbrace{\mathbf{D}_{\omega}^{\ell} \mathbf{A}_{\omega}^{L \rightarrow \ell+1} \boldsymbol{\Sigma}_{\mathbf{x}}^{-1} r_{\omega}^{\ell}(\mathbf{x})}_{\text{residual back-propagated to layer } \ell} \right), \quad (23)$$

$$\text{vect}(\mathbf{W}^{\ell*}) = U_{\omega}^{-1} \text{vect} \left(\underbrace{\sum_{\mathbf{x}} \sum_{\omega} \mathbf{D}_{\omega}^{\ell} \mathbf{A}_{\omega}^{L \rightarrow \ell+1} \boldsymbol{\Sigma}_{\mathbf{x}}^{-1} \left(\mathbf{x} - \sum_{i=\ell}^L \mathbf{A}_{\omega}^{i+1 \rightarrow L} \mathbf{D}_{\omega}^i \mathbf{v}^i \right)}_{\text{residual back-propagated to layer } \ell} \right) (\hat{\mathbf{z}}_{\omega}^{\ell}(\mathbf{x}))^T, \quad (24)$$

we provide detailed derivations below.

G.2.1 Update of the bias parameter

Recall from (2) that $\mathbf{b}_{\omega} = \mathbf{v}^L + \sum_{i=1}^{L-1} \mathbf{W}^L \mathbf{D}_{\omega}^{L-1} \mathbf{W}^{L-1} \dots \mathbf{D}_{\omega}^i \mathbf{v}^i$, we can thus rewrite the loss as

$$\begin{aligned} L(\mathbf{v}^{\ell}) &= -\frac{1}{2} \log \left((2\pi)^{S+D} |\det(\boldsymbol{\Sigma}_{\mathbf{x}})| |\det(\boldsymbol{\Sigma}_{\mathbf{z}})| \right) - \frac{1}{2} \left(\mathbf{x}^T \boldsymbol{\Sigma}_{\mathbf{x}}^{-1} \mathbf{x} - 2 \mathbf{x}^T \boldsymbol{\Sigma}_{\mathbf{x}}^{-1} \left(\sum_{\omega} \mathbf{A}_{\omega} \mathbf{m}_{\omega}^1(\mathbf{x}) + \mathbf{b}_{\omega} \mathbf{m}_{\omega}^0(\mathbf{x}) \right) \right. \\ &\quad \left. + \sum_{\omega} \mathbf{m}_{\omega}^0 \mathbf{b}_{\omega}^T \boldsymbol{\Sigma}_{\mathbf{x}}^{-1} \mathbf{b}_{\omega} + \text{Tr}(\mathbf{A}_{\omega}^T \boldsymbol{\Sigma}_{\mathbf{x}}^{-1} \mathbf{A}_{\omega} \mathbf{M}_{\omega}^2(\mathbf{x})) + 2(\mathbf{A}_{\omega} \mathbf{m}_{\omega}^1(\mathbf{x}))^T \boldsymbol{\Sigma}_{\mathbf{x}}^{-1} \mathbf{b}_{\omega} \right) - \frac{1}{2} \text{Tr}(\boldsymbol{\Sigma}_{\mathbf{z}}^{-1} \mathbf{M}^2(\mathbf{x})) \\ &= -\frac{1}{2} \left(-2 \mathbf{x}^T \boldsymbol{\Sigma}_{\mathbf{x}}^{-1} \left(\sum_{\omega} \mathbf{b}_{\omega} e_{\omega}^0(\mathbf{x}) \right) + \sum_{\omega} e_{\omega}^0 \mathbf{b}_{\omega}^T \boldsymbol{\Sigma}_{\mathbf{x}}^{-1} \mathbf{b}_{\omega} + 2 \sum_{\omega} (\mathbf{A}_{\omega} \mathbf{m}_{\omega}^1(\mathbf{x}))^T \boldsymbol{\Sigma}_{\mathbf{x}}^{-1} \mathbf{b}_{\omega} \right) + cst \\ &= -\frac{1}{2} \sum_{\omega} \left(-2 \mathbf{x}^T \boldsymbol{\Sigma}_{\mathbf{x}}^{-1} (\mathbf{A}_{\omega}^{\ell+1 \rightarrow L} \mathbf{D}_{\omega}^{\ell} \mathbf{v}^{\ell} e_{\omega}^0(\mathbf{x})) + e_{\omega}^0 (\mathbf{A}_{\omega}^{\ell+1 \rightarrow L} \mathbf{D}_{\omega}^{\ell} \mathbf{v}^{\ell})^T \boldsymbol{\Sigma}_{\mathbf{x}}^{-1} (\mathbf{A}_{\omega}^{\ell+1 \rightarrow L} \mathbf{D}_{\omega}^{\ell} \mathbf{v}^{\ell}) \right. \\ &\quad \left. + 2 e_{\omega}^0(\mathbf{x}) \left(\sum_{i \neq \ell} \mathbf{A}_{\omega}^{i+1 \rightarrow L} \mathbf{D}_{\omega}^i \mathbf{v}^i \right)^T \boldsymbol{\Sigma}_{\mathbf{x}}^{-1} (\mathbf{A}_{\omega}^{\ell+1 \rightarrow L} \mathbf{D}_{\omega}^{\ell} \mathbf{v}^{\ell}) + 2 ((\mathbf{m}_{\omega}^1(\mathbf{x}))^T (\mathbf{A}_{\omega})^T \boldsymbol{\Sigma}_{\mathbf{x}}^{-1} \mathbf{A}_{\omega}^{\ell+1 \rightarrow L} \mathbf{D}_{\omega}^{\ell} \mathbf{v}^{\ell}) \right) + cst \\ &= -\frac{1}{2} \sum_{\omega} \left(e_{\omega}^0 (\mathbf{A}_{\omega}^{\ell+1 \rightarrow L} \mathbf{D}_{\omega}^{\ell} \mathbf{v}^{\ell})^T \boldsymbol{\Sigma}_{\mathbf{x}}^{-1} (\mathbf{A}_{\omega}^{\ell+1 \rightarrow L} \mathbf{D}_{\omega}^{\ell} \mathbf{v}^{\ell}) \right) \end{aligned} \quad (A)$$

$$+ 2 (e_{\omega}^0(\mathbf{x}) \left(\sum_{i \neq \ell} \mathbf{A}_{\omega}^{i+1 \rightarrow L} \mathbf{D}_{\omega}^i \mathbf{v}^i - \mathbf{x} \right) + \mathbf{A}_{\omega} e_{\omega}^1(\mathbf{x}))^T \boldsymbol{\Sigma}_{\mathbf{x}}^{-1} (\mathbf{A}_{\omega}^{\ell+1 \rightarrow L} \mathbf{D}_{\omega}^{\ell} \mathbf{v}^{\ell}) \Big) + cst \quad (B)$$

$$\begin{aligned} \Rightarrow \partial L(\mathbf{v}^{\ell}) &= -\frac{1}{2} \sum_{\omega} \left[-e_{\omega}^0(\mathbf{x}) 2 \mathbf{D}_{\omega}^{\ell} \mathbf{A}_{\omega}^{L \rightarrow \ell+1} \boldsymbol{\Sigma}_{\mathbf{x}}^{-1} \mathbf{A}_{\omega}^{\ell+1 \rightarrow L} \mathbf{D}_{\omega}^{\ell} \mathbf{v}^{\ell} \right. \\ &\quad \left. + 2 (\mathbf{A}_{\omega}^{\ell+1 \rightarrow L} \mathbf{D}_{\omega}^{\ell})^T \boldsymbol{\Sigma}_{\mathbf{x}}^{-1} \left(e_{\omega}^0(\mathbf{x}) \left(\sum_{i \neq \ell} \mathbf{A}_{\omega}^{i+1 \rightarrow L} \mathbf{D}_{\omega}^i \mathbf{v}^i - \mathbf{x} \right) + \mathbf{A}_{\omega} e_{\omega}^1(\mathbf{x}) \right) \right] \\ \Rightarrow \mathbf{v}^{\ell} &= \left(\sum_{\mathbf{x}} \sum_{\omega} e_{\omega}^0(\mathbf{x}) \mathbf{D}_{\omega}^{\ell} \mathbf{A}_{\omega}^{L \rightarrow \ell+1} \boldsymbol{\Sigma}_{\mathbf{x}}^{-1} \mathbf{A}_{\omega}^{\ell+1 \rightarrow L} \mathbf{D}_{\omega}^{\ell} \right)^{-1} \\ &\quad \times \sum_{\mathbf{x}} \sum_{\omega \in \Omega} \mathbf{D}_{\omega}^{\ell} \mathbf{A}_{\omega}^{L \rightarrow \ell+1} \boldsymbol{\Sigma}_{\mathbf{x}}^{-1} \left(\mathbf{x} e_{\omega}^0(\mathbf{x}) - \left(\mathbf{A}_{\omega} \mathbf{m}_{\omega}^1(\mathbf{x}) + \sum_{i \neq \ell} \mathbf{m}_{\omega}^0(\mathbf{x}) \mathbf{A}_{\omega}^{i+1 \rightarrow L} \mathbf{D}_{\omega}^i \mathbf{v}^i \right) \right) \end{aligned}$$

as

$$\begin{aligned}
(A) &= e_\omega^0(\mathbf{x})(\mathbf{A}_\omega^{\ell+1 \rightarrow L} \mathbf{D}_\omega^\ell \mathbf{v}^\ell)^T \boldsymbol{\Sigma}_\mathbf{x}^{-1} (\mathbf{A}_\omega^{\ell+1 \rightarrow L} \mathbf{D}_\omega^\ell \mathbf{v}^\ell) \\
\implies \partial(A) &= e_\omega^0(\mathbf{x}) 2 \mathbf{D}_\omega^\ell \mathbf{A}_\omega^{L \rightarrow \ell+1} \boldsymbol{\Sigma}_\mathbf{x}^{-1} \mathbf{A}_\omega^{\ell+1 \rightarrow L} \mathbf{D}_\omega^\ell \mathbf{v}^\ell \\
(B) &= 2 \left[\left(e_\omega^0(\mathbf{x}) \left(\sum_{i \neq \ell} \mathbf{A}_\omega^{i+1 \rightarrow L} \mathbf{D}_\omega^i \mathbf{v}^i - \mathbf{x} \right) + \mathbf{A}_\omega e_\omega^1(\mathbf{x}) \right)^T \boldsymbol{\Sigma}_\mathbf{x}^{-1} (\mathbf{A}_\omega^{\ell+1 \rightarrow L} \mathbf{D}_\omega^\ell \mathbf{v}^\ell) \right] + cst \\
\implies \partial(B) &= (\mathbf{A}_\omega^{\ell+1 \rightarrow L} \mathbf{D}_\omega^\ell)^T \boldsymbol{\Sigma}_\mathbf{x}^{-1} \left(e_\omega^0(\mathbf{x}) \left(\sum_{i \neq \ell} \mathbf{A}_\omega^{i+1 \rightarrow L} \mathbf{D}_\omega^i \mathbf{v}^i - \mathbf{x} \right) + \mathbf{A}_\omega e_\omega^1(\mathbf{x}) \right)
\end{aligned}$$

G.2.2 Update of the slope parameter

We can thus rewrite the loss as

$$\begin{aligned}
L(\mathbf{v}^\ell) &= -\frac{1}{2} \log \left((2\pi)^{S+D} |\det(\boldsymbol{\Sigma}_\mathbf{x})| |\det(\boldsymbol{\Sigma}_\mathbf{z})| \right) - \frac{1}{2} \left(\mathbf{x}^T \boldsymbol{\Sigma}_\mathbf{x}^{-1} \mathbf{x} - 2 \mathbf{x}^T \boldsymbol{\Sigma}_\mathbf{x}^{-1} \left(\sum_\omega \mathbf{A}_\omega m_\omega^1(\mathbf{x}) + \mathbf{b}_\omega m_\omega^0(\mathbf{x}) \right) \right. \\
&\quad \left. + \sum_\omega m_\omega^0 \mathbf{b}_\omega^T \boldsymbol{\Sigma}_\mathbf{x}^{-1} \mathbf{b}_\omega + \text{Tr}(\mathbf{A}_\omega^T \boldsymbol{\Sigma}_\mathbf{x}^{-1} \mathbf{A}_\omega \mathbf{M}_\omega^2(\mathbf{x})) + 2(\mathbf{A}_\omega m_\omega^1(\mathbf{x}))^T \boldsymbol{\Sigma}_\mathbf{x}^{-1} \mathbf{b}_\omega \right) - \frac{1}{2} \text{Tr}(\boldsymbol{\Sigma}_\mathbf{z}^{-1} \mathbf{M}^2(\mathbf{x})) \\
&= \mathbf{x}^T \boldsymbol{\Sigma}_\mathbf{x}^{-1} \left(\sum_\omega \mathbf{A}_\omega e_\omega^1(\mathbf{x}) + \mathbf{b}_\omega e_\omega^0(\mathbf{x}) \right) - \frac{1}{2} \sum_\omega e_\omega^0 \mathbf{b}_\omega^T \boldsymbol{\Sigma}_\mathbf{x}^{-1} \mathbf{b}_\omega - \frac{1}{2} \sum_\omega \text{Tr}(\mathbf{A}_\omega^T \boldsymbol{\Sigma}_\mathbf{x}^{-1} \mathbf{A}_\omega \mathbf{E}_\omega^2(\mathbf{x})) \\
&\quad - \sum_\omega (\mathbf{A}_\omega m_\omega^1(\mathbf{x}))^T \boldsymbol{\Sigma}_\mathbf{x}^{-1} \mathbf{b}_\omega
\end{aligned}$$

Notice that we can rewrite $\mathbf{b}_\omega = \mathbf{A}^{\ell+1 \rightarrow L} \mathbf{D}_\omega^\ell \mathbf{W}^\ell \mathbf{D}_\omega^{\ell-1} \mathbf{b}_\omega^{1 \rightarrow \ell-1} + \sum_{i=\ell}^L \mathbf{A}_\omega^{i+1 \rightarrow L} \mathbf{D}_\omega^i \mathbf{v}^i$ and $\mathbf{A}_\omega = \mathbf{A}^{\ell+1 \rightarrow L} \mathbf{D}_\omega^\ell \mathbf{W}^\ell \mathbf{D}_\omega^{\ell-1} \mathbf{A}^{1 \rightarrow \ell-1}$ and thus we obtain:

$$\begin{aligned}
L(\mathbf{v}^\ell) &= \sum_\omega \mathbf{x}^T \boldsymbol{\Sigma}_\mathbf{x}^{-1} \mathbf{A}^{\ell+1 \rightarrow L} \mathbf{D}_\omega^\ell \mathbf{W}^\ell \mathbf{D}_\omega^{\ell-1} \left(\mathbf{A}^{1 \rightarrow \ell-1} e_\omega^1(\mathbf{x}) + \mathbf{b}_\omega^{1 \rightarrow \ell-1} e_\omega^0(\mathbf{x}) \right) \\
&\quad - \frac{1}{2} \sum_\omega e_\omega^0 (\mathbf{A}^{\ell+1 \rightarrow L} \mathbf{D}_\omega^\ell \mathbf{W}^\ell \mathbf{D}_\omega^{\ell-1} \mathbf{b}_\omega^{1 \rightarrow \ell-1})^T \boldsymbol{\Sigma}_\mathbf{x}^{-1} (\mathbf{A}^{\ell+1 \rightarrow L} \mathbf{D}_\omega^\ell \mathbf{W}^\ell \mathbf{D}_\omega^{\ell-1} \mathbf{b}_\omega^{1 \rightarrow \ell-1}) \\
&\quad - \sum_\omega e_\omega^0 (\mathbf{A}^{\ell+1 \rightarrow L} \mathbf{D}_\omega^\ell \mathbf{W}^\ell \mathbf{D}_\omega^{\ell-1} \mathbf{b}_\omega^{1 \rightarrow \ell-1})^T \boldsymbol{\Sigma}_\mathbf{x}^{-1} \left(\sum_{i=\ell}^L \mathbf{A}_\omega^{i+1 \rightarrow L} \mathbf{D}_\omega^i \mathbf{v}^i \right) \\
&\quad - \frac{1}{2} \sum_\omega \text{Tr} \left((\mathbf{A}_\omega^{\ell+1 \rightarrow L} \mathbf{D}_\omega^\ell \mathbf{W}^\ell \mathbf{D}_\omega^{\ell-1} \mathbf{A}^{1 \rightarrow \ell-1})^T \boldsymbol{\Sigma}_\mathbf{x}^{-1} (\mathbf{A}_\omega^{\ell+1 \rightarrow L} \mathbf{D}_\omega^\ell \mathbf{W}^\ell \mathbf{D}_\omega^{\ell-1} \mathbf{A}^{1 \rightarrow \ell-1}) \mathbf{E}_\omega^2(\mathbf{x}) \right) \\
&\quad - \sum_\omega (\mathbf{A}_\omega^{\ell+1 \rightarrow L} \mathbf{D}_\omega^\ell \mathbf{W}^\ell \mathbf{D}_\omega^{\ell-1} \mathbf{A}^{1 \rightarrow \ell-1} m_\omega^1(\mathbf{x}))^T \boldsymbol{\Sigma}_\mathbf{x}^{-1} (\mathbf{A}_\omega^{\ell+1 \rightarrow L} \mathbf{D}_\omega^\ell \mathbf{W}^\ell \mathbf{D}_\omega^{\ell-1} \mathbf{b}_\omega^{1 \rightarrow \ell-1}) \\
&\quad - \sum_\omega (\mathbf{A}_\omega^{\ell+1 \rightarrow L} \mathbf{D}_\omega^\ell \mathbf{W}^\ell \mathbf{D}_\omega^{\ell-1} \mathbf{A}^{1 \rightarrow \ell-1} m_\omega^1(\mathbf{x}))^T \boldsymbol{\Sigma}_\mathbf{x}^{-1} \left(\sum_{i=\ell}^L \mathbf{A}_\omega^{i+1 \rightarrow L} \mathbf{D}_\omega^i \mathbf{v}^i \right) + cst \\
&= \sum_\omega \mathbf{x}^T \boldsymbol{\Sigma}_\mathbf{x}^{-1} \mathbf{A}^{\ell+1 \rightarrow L} \mathbf{D}_\omega^\ell \mathbf{W}^\ell \mathbf{D}_\omega^{\ell-1} \left(\mathbf{A}^{1 \rightarrow \ell-1} e_\omega^1(\mathbf{x}) + \mathbf{b}_\omega^{1 \rightarrow \ell-1} e_\omega^0(\mathbf{x}) \right) \tag{A}
\end{aligned}$$

$$\begin{aligned}
E &= - \sum_{\omega} (\mathbf{A}_{\omega}^{\ell+1 \rightarrow L} \mathbf{D}_{\omega}^{\ell} \mathbf{W}^{\ell} \mathbf{D}_{\omega}^{\ell-1} \mathbf{A}^{1 \rightarrow \ell-1} m_{\omega}^1(\mathbf{x}))^T \Sigma_{\mathbf{x}}^{-1} (\mathbf{A}^{\ell+1 \rightarrow L} \mathbf{D}_{\omega}^{\ell} \mathbf{W}^{\ell} \mathbf{D}_{\omega}^{\ell-1} \mathbf{b}_{\omega}^{1 \rightarrow \ell-1}) \\
&= - \sum_{\omega} \text{Tr}((\mathbf{W}^{\ell})^T (\mathbf{A}_{\omega}^{\ell+1 \rightarrow L} \mathbf{D}_{\omega}^{\ell})^T \Sigma_{\mathbf{x}}^{-1} \mathbf{A}^{\ell+1 \rightarrow L} \mathbf{D}_{\omega}^{\ell} \mathbf{W}^{\ell} \mathbf{D}_{\omega}^{\ell-1} \mathbf{b}_{\omega}^{1 \rightarrow \ell-1} (\mathbf{D}_{\omega}^{\ell-1} \mathbf{A}^{1 \rightarrow \ell-1} m_{\omega}^1(\mathbf{x}))^T) \\
\Rightarrow \partial E &= - \sum_{\omega} \mathbf{D}_{\omega}^{\ell} \mathbf{A}^{L \rightarrow \ell+1} \Sigma_{\mathbf{x}}^{-1} \mathbf{A}_{\omega}^{\ell+1 \rightarrow L} \mathbf{D}_{\omega}^{\ell} \mathbf{W}^{\ell} (\mathbf{D}_{\omega}^{\ell-1} (\mathbf{b}_{\omega}^{1 \rightarrow \ell-1} (m_{\omega}^1(\mathbf{x}))^T \mathbf{A}_{\omega}^{\ell-1 \rightarrow 1} + \mathbf{A}_{\omega}^{1 \rightarrow \ell-1} m_{\omega}^1(\mathbf{x}) (\mathbf{b}_{\omega}^{1 \rightarrow \ell-1})^T) (\mathbf{D}_{\omega}^{\ell-1})^T)
\end{aligned}$$

we can group B,D and E together as well as A and C. Now to solve this equal 0 we will need to consider the flatten version of \mathbf{W}^{ℓ} which we denote by $\mathbf{w}^{\ell} = \text{vect}(\mathbf{W}^{\ell})$ leading to

$$\begin{aligned}
\partial L &= \sum_{\omega} \mathbf{D}_{\omega}^{\ell} \mathbf{A}_{\omega}^{L \rightarrow \ell+1} \Sigma_{\mathbf{x}}^{-1} (\mathbf{x} - \sum_{i=\ell}^L \mathbf{A}_{\omega}^{i+1 \rightarrow L} \mathbf{D}_{\omega}^i \mathbf{v}^i) (\mathbf{D}_{\omega}^{\ell-1} (\mathbf{A}^{1 \rightarrow \ell-1} \mathbf{e}_{\omega}^1(\mathbf{x}) + \mathbf{b}^{1 \rightarrow \ell-1} \mathbf{e}_{\omega}^0(\mathbf{x})))^T \\
&\quad - \sum_{\omega} \mathbf{D}_{\omega}^{\ell} \mathbf{A}^{L \rightarrow \ell+1} \Sigma_{\mathbf{x}}^{-1} \mathbf{A}_{\omega}^{\ell+1 \rightarrow L} \mathbf{D}_{\omega}^{\ell} \mathbf{W}^{\ell} \mathbf{D}_{\omega}^{\ell-1} \left(\mathbf{e}_{\omega}^0(\mathbf{x}) \mathbf{b}_{\omega}^{1 \rightarrow \ell-1} (\mathbf{b}_{\omega}^{1 \rightarrow \ell-1})^T + \mathbf{A}_{\omega}^{1 \rightarrow \ell-1} \mathbf{E}_{\omega}^2(\mathbf{x}) \mathbf{A}_{\omega}^{\ell-1 \rightarrow 1} \right. \\
&\quad \quad \quad \left. + \mathbf{b}_{\omega}^{1 \rightarrow \ell-1} (m_{\omega}^1(\mathbf{x}))^T \mathbf{A}_{\omega}^{\ell-1 \rightarrow 1} + \mathbf{A}_{\omega}^{1 \rightarrow \ell-1} m_{\omega}^1(\mathbf{x}) (\mathbf{b}_{\omega}^{1 \rightarrow \ell-1})^T \right) \mathbf{D}_{\omega}^{\ell-1} \\
&= \sum_{\omega} P_{\omega}(\mathbf{x})^{\ell} - U_{\omega}^{\ell} \mathbf{W}^{\ell} V_{\omega}^{\ell}(\mathbf{x}) \\
\Rightarrow & \left(\sum_{\mathbf{x}} \sum_{\omega} U_{\omega}^{\ell} \otimes (V_{\omega}^{\ell}(\mathbf{x}))^T \right) \text{vect}(\mathbf{W}^{\ell}) = \sum_{\mathbf{x}} \sum_{\omega} \text{vect}(P_{\omega}^{\ell}(\mathbf{x})) \\
\Rightarrow \text{vect}(\mathbf{W}^{\ell})^* &= \left(\sum_{\mathbf{x}} \sum_{\omega} U_{\omega}^{\ell} \otimes (V_{\omega}^{\ell}(\mathbf{x}))^T \right)^{-1} \left(\sum_{\mathbf{x}} \sum_{\omega} \text{vect}(P_{\omega}^{\ell}(\mathbf{x})) \right)
\end{aligned}$$

H Regularization

We propose in this section a brief discussion on the impact of using a probabilistic prior on the weights of the GDN. In particular, it is clear that imposing a Gaussian prior with zero mean and isotropic covariance on the weights falls back in the log likelihood to impose a $l2$ regularization of the weights with parameter based on the covariance of the prior. If the prior is a Laplace distribution, the log-likelihood will turn the prior into an $l1$ regularization of the weights, again with regularization coefficient based on the prior covariance. Finally, in the case of uniform prior with finite support, the log likelihood will be equivalent to a weight clipping, a standard technique employed in DNs where the weights can not take values outside of a predefined range.

I Computational Complexity

The computational complexity of the method increases drastically with the latent space dimension, and the number of regions, and the number of faces per regions. Those last quantities are directly tied into the complexity (depth and width) of the DGNs. This complexity bottleneck comes from the need to search for all regions, and the need to decompose each region into simplices. As such, the EM learning is not yet suitable for large scale application, however based on the obtained analytical forms, it is possible to derive an approximation of the true form that would be more tractable while providing approximation error bounds as opposed to current methods.

J Additional Experiments

In this section we propose to complement the toy circle experiment from the main paper first we an additional $2d$ case and then with the MNIST dataset.

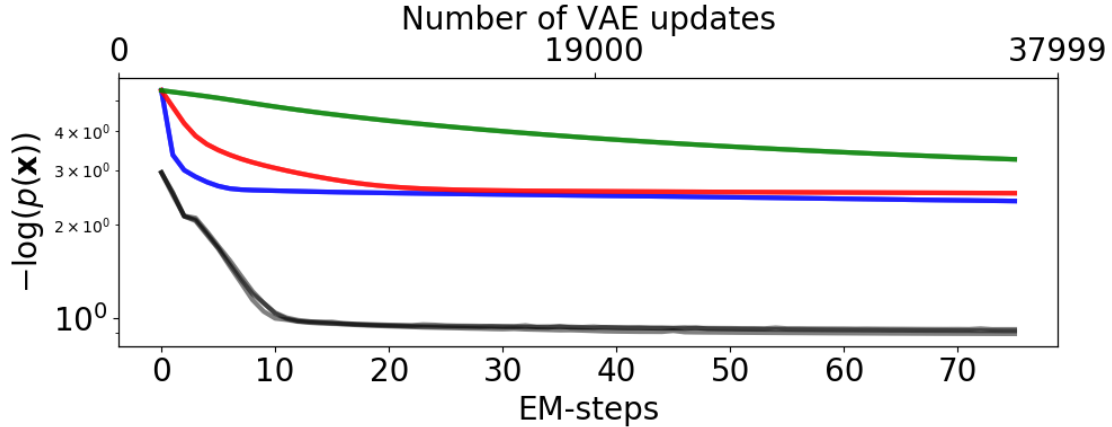


Figure 10: Depiction of the evolution of the NLL during training for the EM and VAE algorithms, we can see that despite the high number of training steps, VAEs are not yet able to correctly approximate the data distribution as opposed to EM training which benefits from much faster convergence. We also see how the VAEs tend to have a large KL divergence between the true posterior and the variational estimate due to this gap, we depict below samples from those models.

Wave

We propose here a simple example where the read data is as follows:

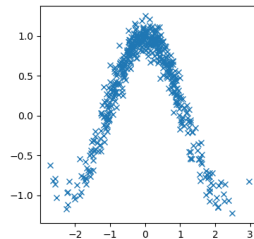


Figure 9: sample of noise data for the wave dataset

We train on this dataset the EM and VAE based learning with various learning rates and depict below the evolution of the NLL for all models, we also depict the samples after learning.

MNIST We now employ MNIST which consists of images of digits, and select the 4 class. Note that due to complexity overhead we maintain a univariate latent space of the GDN and employ a three layer DGN with 8 and 16 hidden units. We provide first the evolution of the NLL through learning for all the training methods and then sample images from the trained DGNs demonstrating how for small DGNs EM learning is able to learn a better data distribution and thus generated realistic samples as opposed to VAEs which need much longer training steps.

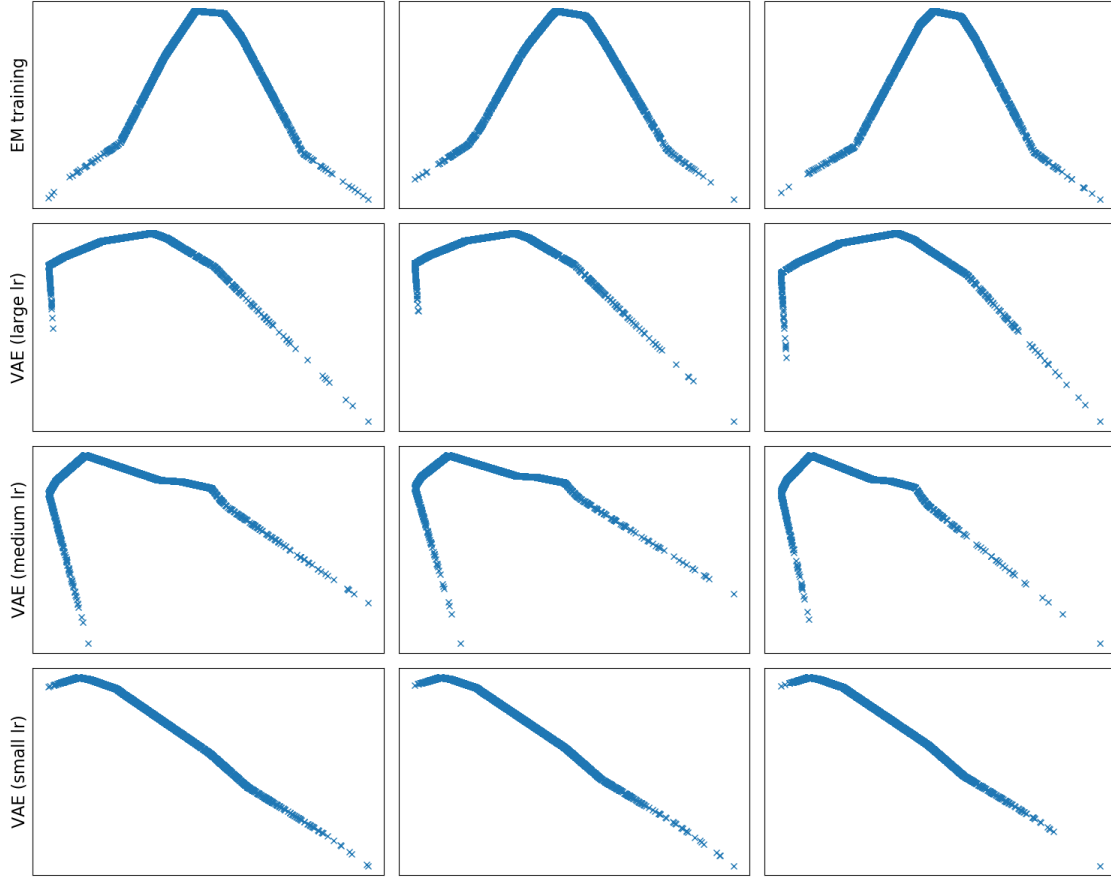


Figure 11: Samples from the various models trained on the wave dataset. We can see on **top** the result of EM training where each column represents a different run, the remaining three rows correspond to the VAE training. Again, EM demonstrates much faster convergence, for VAE to reach the actual data distribution, much more updates are needed.

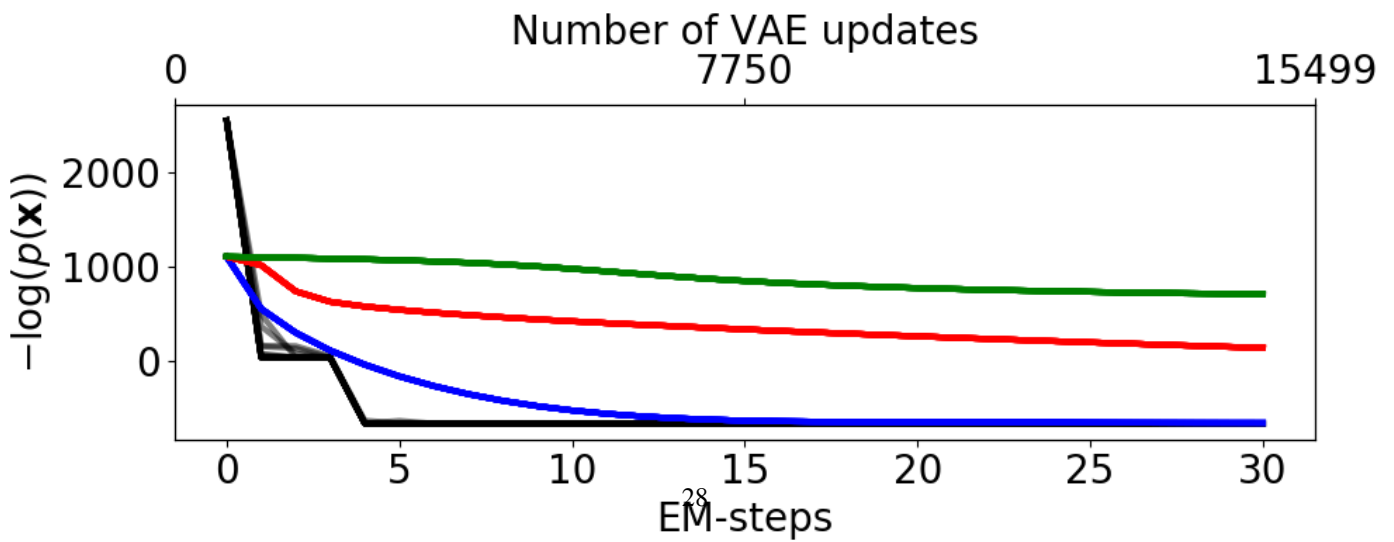


Figure 12: Evolution of the true data negative log-likelihood (in semilog-y plot on MNIST (class 4) for EM and VAE training for a small DGN as described above. The experiments are repeated multiple times, we can see how the learning rate is clearly impacting the learning significantly despite the use of Adam, and that even with the large learning rate, the EM learning is able to reach lower NLL, in fact the quality of the generated samples of the EM modes is much higher as shows below.

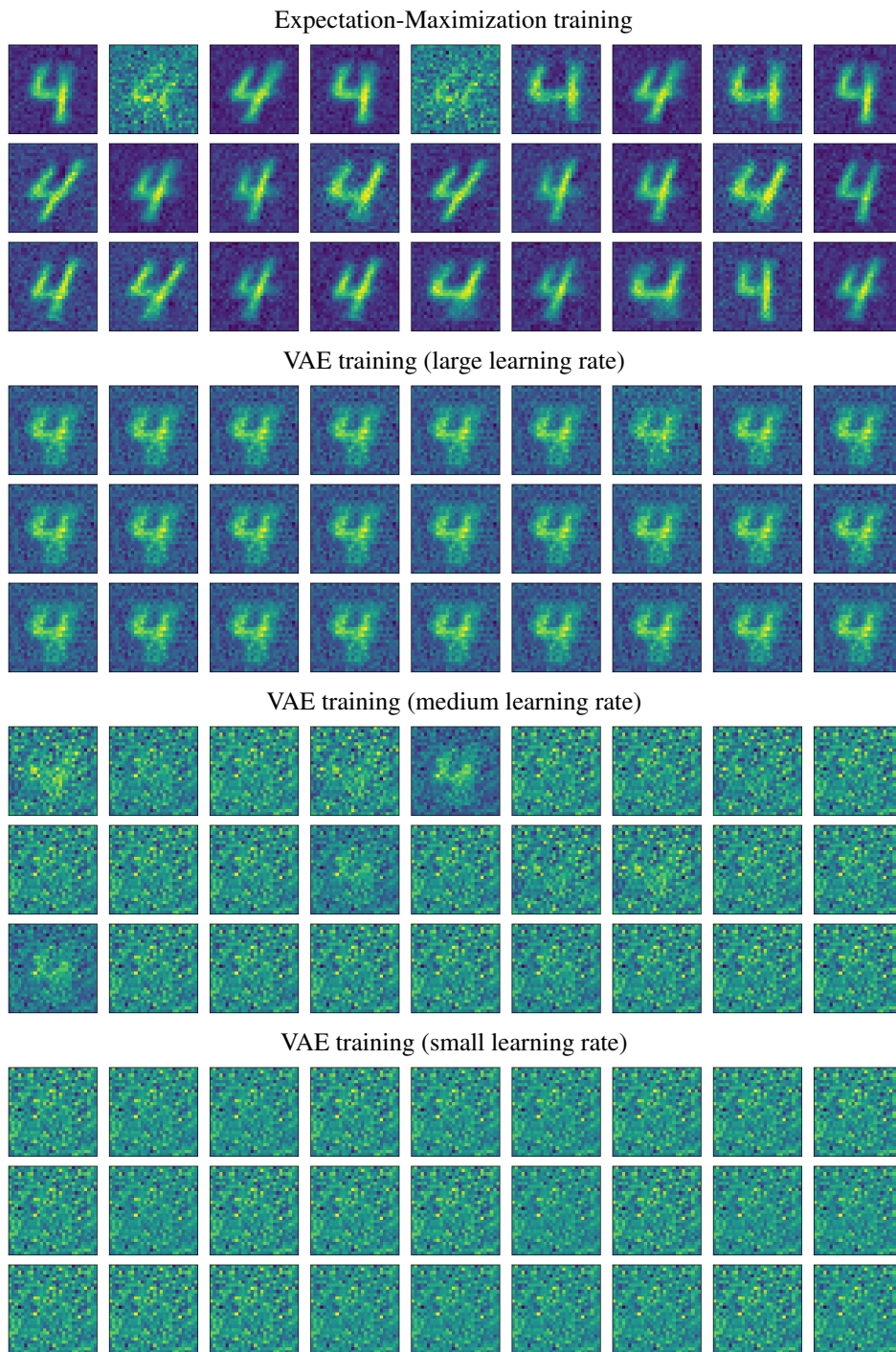


Figure 13: Random samples from trained DGNs with EM or VAEs on a MNIST experiment (with digit 4). We see the ability of EM training to produce realistic and diversified samples despite using a latent space dimension of 1 and a small generative network.