
Mixture-of-experts VAEs can disregard variation in surjective multimodal data

Jannik Wolff^{*†}
TU Berlin

Tassilo Klein, Moin Nabi
SAP AI Research

Rahul G. Krishnan[‡]
University of Toronto

Shinichi Nakajima
TU Berlin

Abstract

Machine learning systems are often deployed in domains that entail data from multiple modalities, for example, phenotypic and genotypic characteristics describe patients in healthcare. Previous works have developed multimodal variational autoencoders (VAEs) that generate several modalities. We consider surjective data, where single datapoints from one modality (such as class labels) describe multiple datapoints from another modality (such as images). We theoretically and empirically demonstrate that multimodal VAEs with a mixture of experts posterior can struggle to capture variability in such surjective data.

1 Introduction

Many datasets entail a surjective mapping between modalities (Fig. 1, “one-to-many data”). That is, an instance from one modality may correspond to several instances from another modality. For example, many computer vision datasets contain labels, attributes, or text data that describe sets of images [LeCun, 1998, Nilsback and Zisserman, 2008, Krizhevsky et al., 2009, Deng et al., 2009, Wah et al., 2011, Liu et al., 2015, Xiao et al., 2017]. Note that “one-to-one data” such as image/caption pairs can become surjective when using data augmentation, e.g., random horizontal flipping of images. Incorporating further modalities can also invoke surjectivity.

Multimodal VAEs maximize a bound on the joint density of several modalities and can thereby learn to generate any modality from any conditioning modality [Suzuki et al., 2016]. For some multimodal VAEs, this bound contains a factor that represents the likelihood of one modality given another modality. We will show that such a factor in the objective function can lead to solutions that disregard heterogeneity within a modality. For example, we demonstrate that samples from models with a mixture of experts posterior such as the MMVAE [Shi et al., 2019] can have a bias towards the class mean of the observed datapoints for a given modality.

2 Method

Let $\mathbf{X} = \{\{\mathbf{x}_m^{(n)}\}_{m=1}^M\}_{n=1}^N$ be a training set with several modalities, where m and n represent the modality and the sample index, respectively. We consider a multimodal VAE with a generative model

$$\begin{aligned} \mathbf{g} &\sim p_\theta(\mathbf{g}), \\ \mathbf{x}_m &\sim p_\theta(\mathbf{x}_m|\mathbf{g}) \quad \text{for } m = 1, \dots, M, \end{aligned} \quad (1)$$

^{*}Correspondence to: wolff.jannik@icloud.com

[†]Part of the work was done at SAP AI Research.

[‡]Part of the work was done at Massachusetts Institute of Technology and Microsoft Research.

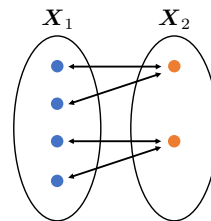


Figure 1: **Surjective data.** X_1 and X_2 depict exemplary modalities. The mapping from the second to the first modality is surjective.

and an inference model

$$\mathbf{g} \sim q_\phi(\mathbf{g}|\{\mathbf{x}_m\}_{m=1}^M). \quad (2)$$

Assume that the generative model (1) is a parametric model, e.g., Gaussian,

$$p_\theta(\mathbf{x}_m|\mathbf{g}) = f_m(\mathbf{x}_m|\boldsymbol{\tau}_m(\mathbf{g};\boldsymbol{\theta})), \quad (3)$$

with the parameters $\{\boldsymbol{\tau}_m\}$, e.g., means and covariances, defined as a function of \mathbf{g} and (typically) neural networks weights $\boldsymbol{\theta}$. Assume that the inference model (2) is defined as a finite mixture with parameters $\boldsymbol{\kappa}_m$ indicating mean and covariance for mixture component r_m (as in the MMVAE [Shi et al., 2019], for example):

$$q_\phi(\mathbf{g}|\{\mathbf{x}_m\}_{m=1}^M) = \frac{1}{M} \sum_{m=1}^M q_\phi(\mathbf{g}|\mathbf{x}_m) = \frac{1}{M} \sum_{m=1}^M r_m(\mathbf{g}|\boldsymbol{\kappa}_m(\mathbf{x}_m; \phi)).$$

Without loss of generality, we assume that \mathbf{x}_M is the label modality, and let $\mathcal{S}_c = \{n \mid \mathbf{x}_M^{(n)} = c\}$ be the set of indices of the samples belonging to the label $c \in \{1, \dots, C\}$. We consider a maximization problem given the following objective function:

$$L_m(\boldsymbol{\theta}, \phi; \mathbf{X}) \equiv \sum_{n=1}^N \int r_M(\mathbf{g}|\boldsymbol{\kappa}_M(\mathbf{x}_M^{(n)}; \phi)) \log f_m(\mathbf{x}_m^{(n)}|\boldsymbol{\tau}_m(\mathbf{g}; \boldsymbol{\theta})) d\mathbf{g}, \quad (4)$$

which is an ELBO for

$$\log p(\mathbf{x}_m|\mathbf{x}_M) = \log \int q_\phi(\mathbf{g}|\mathbf{x}_M) p_\theta(\mathbf{x}_m|\mathbf{g}) d\mathbf{g} \geq \int q_\phi(\mathbf{g}|\mathbf{x}_M) \log p_\theta(\mathbf{x}_m|\mathbf{g}) d\mathbf{g} = L_m(\boldsymbol{\theta}, \phi; \mathbf{X}).$$

Importantly, the MMVAE [Shi et al., 2019] relies on term (4) for learning data translation ability from \mathbf{x}_M to \mathbf{x}_m . Specifically, the authors used stratified sampling for training⁴, which implies that Eq. 4 and term ① from Eq. 5 are related:

$$\begin{aligned} \log p_\theta(\{\mathbf{x}_m\}_{m=1}^M) &\geq \frac{1}{M} \sum_{m=1}^M \mathbb{E}_{q_\phi(\mathbf{g}|\mathbf{x}_m)} \left[\log \frac{p_\theta(\mathbf{g}, \{\mathbf{x}_m\}_{m=1}^M)}{q_\phi(\mathbf{g}|\{\mathbf{x}_m\}_{m=1}^M)} \right] \\ &= \frac{1}{M} \left(\sum_{m=1}^{M-1} \left(\mathbb{E}_{q_\phi(\mathbf{g}|\mathbf{x}_m)} \left[\log \frac{p_\theta(\mathbf{g}, \{\mathbf{x}_m\}_{m=1}^M)}{q_\phi(\mathbf{g}|\{\mathbf{x}_m\}_{m=1}^M)} \right] \right) + \mathbb{E}_{q_\phi(\mathbf{g}|\mathbf{x}_M)} \left[\log \frac{p_\theta(\mathbf{g})}{q_\phi(\mathbf{g}|\{\mathbf{x}_m\}_{m=1}^M)} \right] \right) \\ &\quad + \underbrace{\sum_{i=1}^M \mathbb{E}_{q_\phi(\mathbf{g}|\mathbf{x}_M)} \left[\log p_\theta(\mathbf{x}_i|\mathbf{g}) \right]}_{\textcircled{1}} \end{aligned} \quad (5)$$

The following theorem holds:

Theorem 1. Assume a training set $X = \{\mathbf{x}_m^{(n)}\}_{n \in \mathcal{S}_c}$ which belong to the same label, i.e., $\mathbf{x}_M^{(n)} = c, \forall n \in \mathcal{S}_c$, and there exists $\hat{\boldsymbol{\theta}}$ such that $\boldsymbol{\tau}_m(\mathbf{g}; \hat{\boldsymbol{\theta}})$ is a constant with respect to \mathbf{g} and the maximum likelihood estimator of the parametric model $f_m(\mathbf{x}_m|\boldsymbol{\tau}_m(\mathbf{g}; \boldsymbol{\theta}))$ for the training data. Then, for any $\boldsymbol{\theta}, \phi$, it holds that

$$L_m(\hat{\boldsymbol{\theta}}, \phi; \mathbf{X}) \geq L_m(\boldsymbol{\theta}, \phi; \mathbf{X}). \quad (6)$$

(Proof) Since we assume that $\mathbf{x}_M^{(n)} = c$ for all $n \in \mathcal{S}_c$, the inferred distribution for \mathbf{g} is the same for all n , i.e., $\tilde{r}_M(\mathbf{g}) = r_M(\mathbf{g}|\boldsymbol{\kappa}_M(\mathbf{x}_M^{(n)}; \phi))$. For any such inference model $\tilde{r}_M(\mathbf{g})$, the objective is

⁴Moving Σ_m into the log in Eq. 5 would imply a tighter bound. However, the model may then weigh the experts differently w.r.t. to their gradients, which can disproportionately favor the representation of single modalities at the expense of learning structure across all modalities.

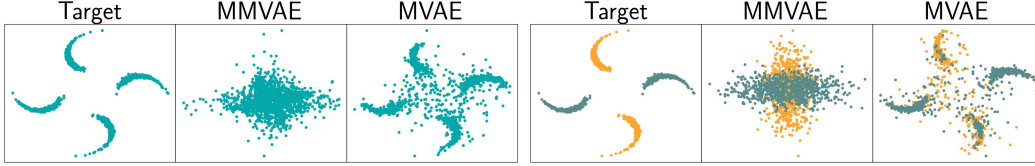


Figure 2: **Generated samples for the first modality.** Left: using samples from $p(\mathbf{g})$. Right: using samples from $q(\mathbf{g}|\mathbf{x}_2)$, where \mathbf{x}_2 are class labels (yellow or green).

upper-bounded by

$$\begin{aligned}
 L_m(\boldsymbol{\theta}, \boldsymbol{\phi}; \mathbf{X}) &= \int \tilde{r}_M(\mathbf{g}; \boldsymbol{\phi}) \left(\sum_{n=1}^N \log f_m(\mathbf{x}_m^{(n)} | \boldsymbol{\tau}_m(\mathbf{g}; \boldsymbol{\theta})) \right) d\mathbf{g} \\
 &\leq \int \tilde{r}_M(\mathbf{g}; \boldsymbol{\phi}) \left(\sum_{n=1}^N \log f_m(\mathbf{x}_m^{(n)} | \hat{\boldsymbol{\tau}}_m) \right) d\mathbf{g}
 \end{aligned} \tag{7}$$

with the maximum likelihood estimator $\hat{\boldsymbol{\tau}}_m$ for the parametric model f_m given the training set $\{\mathbf{x}_m^{(n)}\}_{n=1}^{S_c}$. The assumed existence of $\hat{\boldsymbol{\theta}}$ such that $\boldsymbol{\tau}_m(\mathbf{g}; \hat{\boldsymbol{\theta}}) = \hat{\boldsymbol{\tau}}_m$ leads to Eq. (6). \square

Intuitively, consider a single class: $c \in \{1\}$. Let $p_\theta(\mathbf{x}_m|\mathbf{g})$ be Gaussian with diagonal covariance, where $\mathbf{g} \sim q_\phi(\mathbf{g}|\mathbf{x}_M)$. Theorem 1 implies the existence of an upper bound where the mean parameter from $p_\theta(\mathbf{x}_m|\mathbf{g})$ always coincides with the mean from $\{\mathbf{x}_m^{(n)}\}_{n \in S_c}$ for any \mathbf{g} . This solution is invariant to \mathbf{g} because \mathbf{x}_M does not carry information about across-datapoint variability in \mathbf{x}_m . In other words, the solution maximizes the likelihood of the training data $\{\mathbf{x}_m^{(n)}\}_{n=1}^{S_c}$ with a single Gaussian distribution. That is, the mean parameter minimizes the distance to all datapoints from modality m simultaneously: the model captures the mean of the target distribution – not its variability.

3 Experiments

We create a synthetic dataset (inspired by Johnson et al. [2016]) with modality $\mathbf{x}_1 \in \mathbb{R}^2$ and label modality $\mathbf{x}_2 \in \{0, 1\}$. We implement the MVAE [Wu and Goodman, 2018] and MMVAE [Shi et al., 2019]. The latent distributions are isotropic Gaussian. The generative distributions are isotropic Gaussian for the first modality and categorical for the second modality.

For the MMVAE, Fig. 2 supports our argument that samples for the first modality tend towards the mean of the observed datapoints (for the same class). The MVAE does not suffer from this problem, possibly because the MVAE’s objective function does not contain the factor $p(\mathbf{x}_1|\mathbf{x}_2)$ (App. A). App. B visualizes the latent spaces, which are two-dimensional to avoid possible obfuscation from dimensionality-reduction techniques.

4 Conclusion

We show that multimodal VAEs with a mixture posterior can struggle to capture heterogeneity in surjective data. This finding implies that practitioners should closely consider the type of data when training such models: for example, data augmentation may not be beneficial since this procedure often promotes surjectivity. Future work may investigate possible solutions, e.g., by considering models that do not maximize $p(\mathbf{x}_m|\mathbf{x}_{M \neq m})$ explicitly. It would be interesting to analyze how such a solution affects robustness.

Acknowledgements

SN is supported by the German Ministry for Education and Research as BIFOLD - Berlin Institute for the Foundations of Learning and Data (ref. 01IS18025A and ref. 01IS18037A). RGK was supported by a grant from SAP Corporation.

References

- J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.
- M. J. Johnson, D. K. Duvenaud, A. Wiltchko, R. P. Adams, and S. R. Datta. Composing graphical models with neural networks for structured representations and fast inference. In *Advances in neural information processing systems*, pages 2946–2954, 2016.
- A. Krizhevsky, G. Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- Y. LeCun. The mnist database of handwritten digits. <http://yann.lecun.com/exdb/mnist/>, 1998.
- Z. Liu, P. Luo, X. Wang, and X. Tang. Deep learning face attributes in the wild. In *Proceedings of International Conference on Computer Vision (ICCV)*, December 2015.
- M.-E. Nilsback and A. Zisserman. Automated flower classification over a large number of classes. In *2008 Sixth Indian Conference on Computer Vision, Graphics & Image Processing*, pages 722–729. IEEE, 2008.
- Y. Shi, N. Siddharth, B. Paige, and P. Torr. Variational mixture-of-experts autoencoders for multimodal deep generative models. In *Advances in Neural Information Processing Systems*, pages 15692–15703, 2019.
- M. Suzuki, K. Nakayama, and Y. Matsuo. Joint multimodal learning with deep generative models. *arXiv preprint arXiv:1611.01891*, 2016.
- C. Wah, S. Branson, P. Welinder, P. Perona, and S. Belongie. The caltech-ucsd birds-200-2011 dataset. 2011.
- M. Wu and N. Goodman. Multimodal generative models for scalable weakly-supervised learning. In *Advances in Neural Information Processing Systems*, pages 5575–5585, 2018.
- H. Xiao, K. Rasul, and R. Vollgraf. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.

A Theorem 1 does not apply to the MVAE

The MVAE [Wu and Goodman, 2018] employs a product posterior inspired by the true posterior:

$$q_\phi(\mathbf{g}|\{\mathbf{x}_m\}_{m=1}^M) \propto p_\theta(\mathbf{g}) \prod_{m=1}^M q_\phi(\mathbf{g}|\mathbf{x}_m). \quad (8)$$

In our experiments from § 3, we follow Wu and Goodman [2018] and maximize the following three ELBOs:

$$L(\theta, \phi; \mathbf{X}) := ELBO(\mathbf{x}_1, \mathbf{x}_2) + ELBO(\mathbf{x}_1) + ELBO(\mathbf{x}_2) \quad (9)$$

The ELBO for M modalities is defined as:

$$\begin{aligned} ELBO(\{\mathbf{x}_m\}_{m=1}^M) &:= \mathbb{E}_{q_\phi(\mathbf{g}|\{\mathbf{x}_m\}_{m=1}^M)} \left[\log \frac{p_\theta(\mathbf{g})}{q_\phi(\mathbf{g}|\{\mathbf{x}_m\}_{m=1}^M)} \right] + \sum_{m=1}^M \mathbb{E}_{q_\phi(\mathbf{g}|\{\mathbf{x}_m\}_{m=1}^M)} [\log p_\theta(\mathbf{x}_m|\mathbf{g})] \\ &\leq \log p_\theta(\{\mathbf{x}_m\}_{m=1}^M), \end{aligned} \quad (10)$$

Therefore, $p_\theta(\mathbf{x}_m|\mathbf{g})$ is always conditioned on \mathbf{x}_m via the importance distribution, i.e., the model learns $p(\mathbf{x}_m|\{\mathbf{x}_i\}_{i=1}^M)$ or $p(\mathbf{x}_m|\mathbf{x}_m)$. This implies that the MVAE does not explicitly optimize $p(\mathbf{x}_{m \neq M}|\mathbf{x}_M)$ for any $m \neq M$, i.e., Theorem 1 does not apply to the MVAE.

B Additional experimental results

The solution $q(\mathbf{g}|\mathbf{x}_1) = q(\mathbf{g}|\mathbf{x}_2)$ can be helpful because it implies that samples from either posterior produce the same generative distribution for any modality. Figure 3 indicates that the MVAE aligns these marginal posteriors better than the MMVAE, which possibly explains the MVAE’s better generative capability in Fig. 2. Figure 2 further exposes that even the MVAE struggles to represent the data perfectly. Its latent representations from Fig. 3 reveal that the model produces some overlap between the class manifolds of the marginal posteriors for the second modality – possibly in an attempt to fit the isotropic Gaussian prior $p(\mathbf{g})$. We assume that this struggle is caused by the fact that there are just two unique label datapoints.

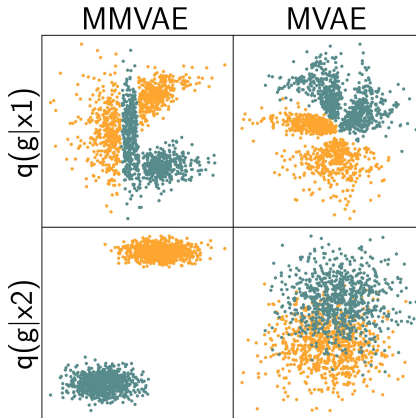


Figure 3: Marginal posteriors over the latent variable \mathbf{g} .