

---

# Information Maximization Auto-Encoding

---

**Dejiao Zhang**

University of Michigan, Ann Arbor  
dejiao@umich.edu

**Tianchen Zhao**

University of Michigan, Ann Arbor  
ericolon@umich.edu

**Laura Balzano**

University of Michigan, Ann Arbor  
girasole@umich.edu

## Abstract

We propose the Information Maximization Autoencoder (IMAE), an information theoretic approach to simultaneously learn continuous and discrete representations in an unsupervised setting. Unlike the Variational Autoencoder framework, IMAE starts from a stochastic encoder that seeks to map each input data to a hybrid discrete and continuous representation with the objective of maximizing the mutual information between the data and their representations. A decoder is included to approximate the posterior distribution of the data given their representations, where a high fidelity approximation can be achieved by leveraging the informative representations. We show that the proposed objective is theoretically valid and provides a new perspective for understanding the tradeoffs regarding informativeness of the representation factors, disentanglement of representations, and decoding quality.

## 1 Introduction

A central tenet for designing and learning a model for data is that the resulting representation should be compact yet informative. Therefore, the goal of learning can be formulated as finding informative representations about the data under proper constraints. In this work, we propose an information theoretic approach to simultaneously learn continuous and discrete representations in an unsupervised setting. We start with a stochastic encoder  $p_\theta(z|x)$  and aim at maximizing the mutual information between the data  $x$  and its representation  $z$ . In this setting, a reconstruction or generating phase can be obtained as the variational inference of the true posterior  $p_\theta(x|z)$ . By explicitly seeking informative representations, the proposed model yields better decoding quality. Moreover, we show that the information maximization objective naturally induces a balance between the informativeness of each latent factor and the statistical independence between them, which gives a more principled way to learn semantically meaningful representations.

Another contribution of this work is proposing a framework for simultaneously learning continuous and discrete representations for categorical data. Categorical data are ubiquitous in real-world tasks, where using a hybrid discrete and continuous representation to capture both categorical information and continuous variation in data is more consistent with the natural generation process. We focus on categorical data that are similar in nature, *i.e.*, where different categories still share similar variations (features). We seek to learn semantically meaningful discrete representations while maintaining disentanglement of the continuous representations that capture the variation shared across categories.

## 2 Information Maximization Representation Learning

Given data  $x \in \mathbb{R}^d$ , we consider learning a hybrid continuous-discrete representation, denoted respectively with variables  $z \in \mathbb{R}^{K_1}$  and  $y \in \{1, \dots, K_2\}$ , using a stochastic encoder parameterized by  $\theta$ , *i.e.*,  $p_\theta(y, z|x)$ . We seek to learn compact yet semantically meaningful representations in the

sense that they should be low dimensional but informative enough about the data. A natural approach is to maximize the mutual information [Cover and Thomas, 2012]  $I_\theta(\mathbf{x}; \mathbf{y}, \mathbf{z})$  between the data and its representation under the constraint  $K_1, K_2 \ll d$ . A probabilistic decoder  $q_\phi(\mathbf{x}|\mathbf{y}, \mathbf{z})$  is adopted to approximate the true posterior  $p_\theta(\mathbf{x}|\mathbf{y}, \mathbf{z})$ , which can be hard to estimate or even intractable. The dissimilarity between them is optimized by minimizing the KL divergence  $D_{\text{KL}}(p_\theta(\mathbf{x}|\mathbf{y}, \mathbf{z})||q_\phi(\mathbf{x}|\mathbf{y}, \mathbf{z}))$ . In summary, IMAE considers the following,

$$\text{maximize}_{\theta, \phi} \beta_0 I_\theta(\mathbf{x}; \mathbf{y}, \mathbf{z}) - D_{\text{KL}}(p_\theta(\mathbf{x}|\mathbf{y}, \mathbf{z})||q_\phi(\mathbf{x}|\mathbf{y}, \mathbf{z})) . \quad (1)$$

Given that  $H(x)$  is independent of the optimization procedure, we can show that optimizing (1) is equivalent to optimizing the following:

$$\text{maximize}_{\theta, \phi} \beta I_\theta(\mathbf{x}; \mathbf{y}, \mathbf{z}) + \mathbb{E}_{p_\theta(\mathbf{x}, \mathbf{y}, \mathbf{z})} [\log q_\phi(\mathbf{x}|\mathbf{y}, \mathbf{z})], \quad \beta = \beta_0 - 1 > 0 . \quad (2)$$

We set  $\beta > 0$  to balance between maximizing the informativeness of latent representations and maintaining the decoding quality. The second term is often interpreted as the ‘‘reconstruction error’’ which can be optimized using the reparameterization tricks proposed by [Kingma and Welling, 2013] and [Jang et al., 2016] for the continuous representation  $\mathbf{z}$  and the discrete representation  $\mathbf{y}$  respectively. Now we introduce our method to optimize the first term  $I_\theta(\mathbf{x}; \mathbf{y}, \mathbf{z})$  in (2).

## 2.1 Simultaneously seeking informativeness and disentanglement

Assuming the conditional distribution of the representation  $(\mathbf{y}, \mathbf{z})$  given  $\mathbf{x}$  is factorial, and similarly for the marginal distribution  $p_\theta(\mathbf{y}, \mathbf{z}) = p_\theta(\mathbf{y})p_\theta(\mathbf{z})$ , then

$$I_\theta(\mathbf{x}; \mathbf{y}, \mathbf{z}) = I_\theta(\mathbf{x}; \mathbf{y}) + \sum_{k=1}^{K_1} I_\theta(\mathbf{x}; \mathbf{z}_k) - D_{\text{KL}}\left(p_\theta(\mathbf{z})||\prod_{k=1}^{K_1} p_\theta(\mathbf{z}_k)\right) . \quad (3)$$

The first two terms of the RHS quantify how much information each latent factor carries about the data. The last term is known as the *total correlation* of  $\mathbf{z}$  [Watanabe, 1960], which quantifies the statistical independence between the continuous latent factors and achieves the minimum if and only if they are independent of each other. Eq (3) implies that maximizing  $I_\theta(\mathbf{x}; \mathbf{y}, \mathbf{z})$  can be conducted by maximizing informativeness of each latent factor while simultaneously promoting statistical independence between the continuous factors. Various Monte Carlo based sampling methods have been proposed to optimize the total correlation term [Chen et al., 2018, Esmaeili et al., 2018]; we follow this line. Next we construct tractable approximations for  $I_\theta(\mathbf{x}; \mathbf{z}_k)$  and  $I_\theta(\mathbf{x}; \mathbf{y})$  respectively.

## 2.2 Informative continuous representations

Without any constraints, the mutual information  $I_\theta(\mathbf{x}; \mathbf{z}_k)$  between the data and its continuous representation factor can be trivially maximized by severely fragmenting and exploding the latent space. Specifically, as shown in Proposition 1<sup>1</sup>,  $\mathbf{z}_k$  is more informative about  $\mathbf{x}$  if it has less uncertainty given  $\mathbf{x}$  yet captures more variance in data, *i.e.*,  $\sigma_k(\mathbf{x})$  is small while  $\mu_k(\mathbf{x})$  is dispersed within a large space. This can result in discontinuity of the latent representation  $\mathbf{z}_k$ , where in the extreme case each data point is associated with a delta distribution in the latent space  $p_\theta(\mathbf{z}_k|\mathbf{x}^{(i)}) = \delta(\mathbf{z}_k^{(i)})$ .

**Proposition 1** *Suppose the conditional distribution  $p_\theta(\mathbf{z}|\mathbf{x})$  is a factorial Gaussian distribution with mean  $\mu(\mathbf{x})$  and covariance  $\Sigma(\mathbf{x})$ . Let  $\sigma(\mathbf{x}) \in \mathbb{R}^{K_1}$  denote the diagonal entries of  $\Sigma(\mathbf{x})$ , then*

$$I_\theta(\mathbf{x}; \mathbf{z}_k) \leq \frac{1}{2} \log \left[ (\mathbb{E}_\mathbf{x} [\sigma_k^2(\mathbf{x})] + \text{Var}_\mathbf{x} [\mu_k(\mathbf{x})]) \right] - \frac{1}{2} \mathbb{E}_\mathbf{x} [\log \sigma_k^2(\mathbf{x})] , \quad k = 1, \dots, K_1 . \quad (4)$$

The equality in (4) is attained if and only if  $\mathbf{z}_k$  is Gaussian distributed, given which we have

$$I_\theta(\mathbf{x}; \mathbf{z}_k) \geq \frac{1}{2} \log (1 + \text{Var}_\mathbf{x} [\mu_k(\mathbf{x})] / \mathbb{E}_\mathbf{x} [\sigma_k^2(\mathbf{x})]) , \quad k = 1, \dots, K_1 . \quad (5)$$

To remedy this issue while achieving the upper bound in Proposition 1, a natural resolution is to squeeze  $\mathbf{z}_k$  within the domain of a Gaussian distribution with finite variance. By doing so, we can get a more reasonable trade-off between enlarging the spread of  $\mu_k(\mathbf{x})$  and maintaining the continuity of  $\mathbf{z}_k$ , while achieving the maximal  $I(\mathbf{x}; \mathbf{z})$  among all possible solutions with the same variance of  $\mathbf{z}_k$ . Therefore, we consider the following as the surrogate for maximizing  $I_\theta(\mathbf{x}; \mathbf{z}_k)$ ,

$$\max_\theta \mathcal{L}_\theta(\mathbf{z}) := -\sum_{k=1}^{K_1} D_{\text{KL}}(p_\theta(\mathbf{z}_k)||r(\mathbf{z}_k)) . \quad (6)$$

Here  $r(\mathbf{z}_k)$  is an i.i.d scaled normal distribution with variance being some prefixed finite value.

<sup>1</sup>While similar results have likely been established, we include Proposition 1 to motivate our objective design.

### 2.3 Informative discrete representations

The mutual information  $I_\theta(\mathbf{x}; \mathbf{y})$  between the data and its discrete representation can be well approximated, given the fact that the cardinality of the space of  $\mathbf{y}$  is typically low. To be more precise, Proposition 2 shows that, with a suitably large batch of samples, the empirical mutual information  $\widehat{I}_\theta(\mathbf{x}; \mathbf{y})$  is a good approximation to  $I_\theta(\mathbf{x}; \mathbf{y})$ . This enables us to optimize  $I_\theta(\mathbf{x}; \mathbf{y})$  in a theoretically justifiable way that is amenable to stochastic gradient descent with minibatches of data. Hence, to maximize  $I_\theta(\mathbf{x}; \mathbf{y})$  we consider the following,

$$\max_{\theta} \mathcal{L}_\theta(\mathbf{y}) := \widehat{I}_\theta(\mathbf{x}; \mathbf{y}) = \mathbf{H} \left( \frac{1}{M} \sum_{m=1}^M p_\theta(\mathbf{y}|x_m) \right) - \frac{1}{M} \sum_{m=1}^M \mathbf{H}(p_\theta(\mathbf{y}|x_m)). \quad (7)$$

**Proposition 2** *Assume the marginal probabilities of the true and predicted labels are bounded below, i.e.  $p_\theta(\mathbf{y}), \widehat{p}_\theta(\mathbf{y}) \in [1/(CK_2), 1]$  for all  $\mathbf{y} \in \{1, \dots, K_2\}$  with some constant  $C > 1$ . Then  $\mathbb{P} \left( \left| I_\theta(\mathbf{x}; \mathbf{y}) - \widehat{I}_\theta(\mathbf{x}; \mathbf{y}) \right| \leq K_2 (\max\{\log CK_2 - 1, 1\} + e) \sqrt{\frac{\log(2K_2/\delta)}{2M}} \right) \geq 1 - 2\delta, \delta \in (0, 1)$ .*

**Overall Objective** As a summary of (3) (6) and (7), so far we have

$$\max_{\theta, \phi} \beta \left( \mathcal{L}_\theta(\mathbf{z}) + \mathcal{L}_\theta(\mathbf{y}) - D_{\text{KL}} \left[ p(\mathbf{z}) \parallel \prod_{k=1}^{K_1} p(\mathbf{z}_k) \right] \right) + \mathbb{E}_{p_\theta(\mathbf{x}, \mathbf{y}, \mathbf{z})} [\log q_\phi(\mathbf{x}|\mathbf{y}, \mathbf{z})].$$

The first three terms are associated with our information maximization objective, while the last one aims at better approximation of the posterior  $p_\theta(\mathbf{x}|\mathbf{y}, \mathbf{z})$ . A better balance between these two targets can be achieved by weighting them differently. The informativeness of each latent factor can be optimized through  $\mathcal{L}_\theta(\mathbf{z})$  and  $\mathcal{L}_\theta(\mathbf{y})$ , while statistically independent latent continuous factors can be promoted by minimizing the total correlation term. Therefore, trade-offs can be formalized regarding the informativeness of each latent factor, disentanglement of the representation, and better decoding quality. With this final adjustment, we settle on the following overall objective:

$$\begin{aligned} \max_{\theta, \phi} \mathcal{L}_{\text{IMAE}} := & \mathbb{E}_{p_\theta(\mathbf{x}, \mathbf{y}, \mathbf{z})} [\log q_\phi(\mathbf{x}|\mathbf{y}, \mathbf{z})] \\ & + \beta \mathcal{L}_\theta(\mathbf{y}) + \beta \mathcal{L}_\theta(\mathbf{z}) - \gamma D_{\text{KL}} \left[ p_\theta(\mathbf{z}) \parallel \prod_{k=1}^{K_1} p_\theta(\mathbf{z}_k) \right], \quad \beta, \gamma > 0. \end{aligned} \quad (8)$$

## 3 Numerical Results

We compare IMAE against a number of VAE based approaches, including  $\beta$ -VAE [Higgins et al., 2016], InfoVAE [Zhao et al., 2017] and JointVAE [Dupont, 2018]. A summarization of these models is provided in Appendix A. We would like to show that IMAE can (i) successfully learn a hybrid continuous and discrete representation, with  $\mathbf{y}$  matching the natural categorical information well and  $\mathbf{z}$  capturing the disentangled feature information shared across categories; (ii) outperform the other models on achieving a better trade-off between the interpretability of the representation and the decoding quality. We choose the priors  $r(\mathbf{z})$  and  $r(\mathbf{y})$  as isotropic Gaussian distribution and uniform distribution over categories respectively.

### 3.1 Informative representations yield better interpretability

Figure 1 validates Proposition 1 by showing that, with roughly same amount of variance for each latent variable  $z_k$ , those achieving high mutual information with the data have mean values  $\mu_k(\mathbf{x})$  of the conditional probability  $p(z_k|\mathbf{x})$  disperse across data samples and variances  $\sigma_k(\mathbf{x})$  decrease to small values for all data samples. As seen in Figure 1(b)-(d), informative variables in the continuous representation have uncovered intuitive latent factors of the variation in the data, while the factor  $z_8$  has no mutual information with the data and shows no variation. We observe the same phenomenon for the discrete representation  $\mathbf{y}$  in Figure 1(e)&(f), which were obtained with two different values of  $\beta$  and  $\gamma$ , where the more informative one matches the natural labels better.

### 3.2 Quantitative comparisons

In this section, we perform quantitative evaluations on MNIST and dSprites. Before we present our main results, we first describe an assumption that we make on the discrete representations. A

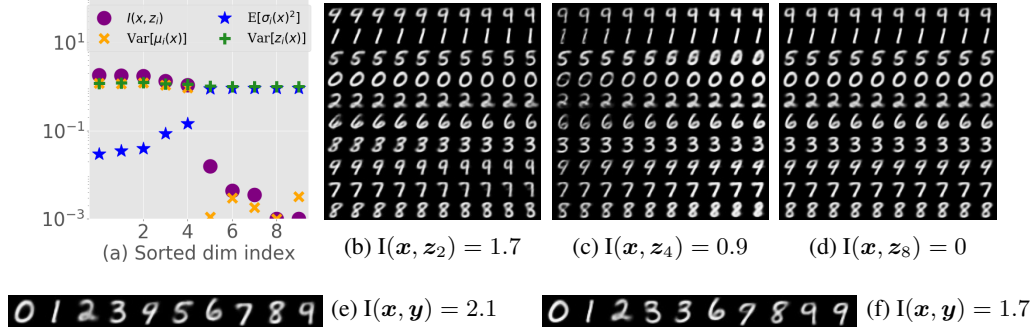


Figure 1: **IMAE on MNIST.** (a) Illustration of Proposition 1. (b)-(d) Latent traverse on the continuous representations  $z$ . The rows are conditioned on the discrete representations  $y$  learnt by IMAE, where for each row  $z$  is initialized by feeding the encoder with randomly selected data whose labels are predicted as  $y$ . We then manipulate each selected  $z_k$  within  $[-2, 2]$  while keeping all other dimensions fixed. (e) & (f) Discrete representations learnt by IMAE with different  $\beta$  values.

reasonable assumption is that the conditional distribution  $p(y|x)$  should be locally smooth so that the data samples that are close on their manifold should have high probability of being assigned to the same category [Agakov, 2005]. This assumption is crucial for using neural networks to learn discrete representations, since it's easy for a high capacity model to learn a non-smooth function  $p(y|x)$  that can abruptly change its predictions without guaranteeing similar data samples will be mapped to similar  $y$ . To remedy this issue, we adopt the virtual adversarial training (VAT) trick proposed by [Miyato et al., 2016] and augment  $\mathcal{L}_\theta(y)$  as follows:<sup>2</sup>

$$\max \mathcal{L}_\theta(y) := \widehat{I}_\theta(x; y) - \mathbb{E}_{\widehat{p}(x)} \left[ \max_{\|\eta\| \leq \epsilon} \mathbf{H}(p_\theta(y|x); p_\theta(y|x + \eta)) \right]. \quad (9)$$

The second term of RHS regularizes  $p_\theta(y|x)$  to be consistent within the  $\epsilon$  norm ball of each data sample so as to maintain local smoothness of the prediction model. *For fair comparison, we augment all four methods with VAT. We found that using VAT is essential for all of them except  $\beta$ -VAE to learn interpretable discrete representations.*

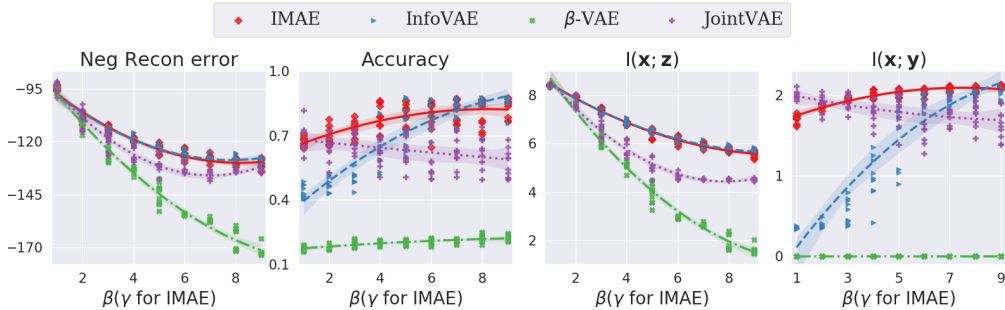


Figure 2: **IMAE on MNIST.** We track the key quantities for different models by sweeping  $\beta$ . We set  $\gamma = 2\beta$  for IMAE. For each  $\beta$ , we run each method over 10 random initializations.

**MNIST** Figure 4 shows that, by simply pushing the conditional distribution  $p(y|x)$  towards the uniform distribution  $r(y)$ ,  $\beta$ -VAE sacrifices the mutual information  $I(x; y)$  and hence struggles in learning interpretable discrete representation even with VAT. On the other hand, large  $\beta$  values drive  $\beta$ -VAE to sacrifice more mutual information  $I(x; z)$  between the data and its continuous representations, which together with the less informative discrete representation result in poor decoding quality. In contrast, the other three methods can remedy this issue to different degrees, and hence attain better trade-off regarding informativeness of latent representations and decoding quality. Compared to JointVAE and InfoVAE, IMAE is more capable of learning discrete representations over a wide range of  $\beta, \gamma$  values, which implies that less overlap between the manifolds of different categories

<sup>2</sup>In this paper, we set  $\epsilon = 1$  across datasets. VAT can be effectively approximated by a pair of forward and backward passes [Miyato et al., 2016].

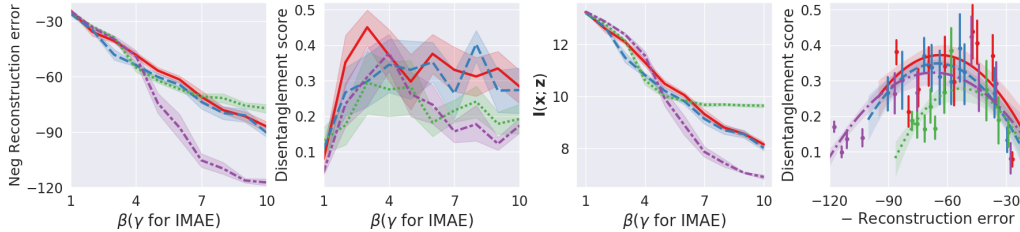


Figure 3: **Disentanglement comparison on dSprites.** IMAE performs well regarding the disentanglement score vs. decoding quality trade-off, especially in the region of interest where both decoding quality and informativeness of representations are fairly good. The results are reported by training each method with  $\beta \in [1, 10]$  (we set  $\beta = \gamma/2$  with  $\gamma \in [1, 10]$  for IMAE). For each  $\beta$  value, every method is trained over 10 random initializations. Shade regions indicate the 80% confidence intervals.

is induced. As a result, IMAE is expected to yield better decoding quality for each category, seen for MNIST in Figure 2. Although InfoVAE and JointVAE also learn comparatively good discrete representations with large and small  $\beta$  values respectively, the corresponding results of these two regions have poor decoding quality or a much lower disentanglement score (see Figure 3). In contrast, IMAE consistently performs well with different hyperparameter values, especially in the region of interest where the decoding quality and the informativeness of latent representations are good enough.

**2D Shapes** We now evaluate IMAE on dSprites where the ground truth of both continuous and discrete representation factors are available. We use the disentanglement metric proposed by [Chen et al., 2018]. The higher the disentanglement score is, the more disentangled the representation factors are.<sup>3</sup>

As shown in Figure 3, with large  $\beta$  values,  $\beta$ -VAE penalizes the mutual information too much and this degrades the usefulness of representations, while all other three methods achieve higher disentanglement score with better decoding quality. For JointVAE, higher  $\beta$  values push the upper bound of the mutual information to converge to the prefixed target value; it therefore can maintain more mutual information between the data and the overall latent representations and give better decoding quality. However, the associated disentanglement quality is poor. This implies that simply restricting the overall capacity of the latent representations is not enough for learning disentangled representations. While InfoVAE yields a comparatively better disentanglement score by pushing the marginal joint distribution of the representations towards a factorial distribution more aggressively with large values of  $\beta$ , the associated decoding quality and informativeness of latent representations are both poor. In contrast, IMAE is capable of achieving a better trade-off between the disentanglement score and the decoding quality in the region of interest where the decoding quality as well as the informativeness are fairly good. We attribute this to the effect of explicitly seeking statistically independent latent factors by minimizing the total correlation term in our objective.

## 4 Conclusion

Unsupervised joint learning of disentangled continuous and discrete representations is a challenging problem due to the lack of a prior for semantic awareness and other inherent difficulties that arise in learning discrete representations. This work takes a step towards achieving this goal. A limitation of our model is that it pursues disentanglement by assuming or trying to encourage independent scalar latent factors, which may not always be sufficient for representing real data. For example, data may exhibit category specific variation, or a subset of latent factors might be correlated. This motivates us to explore more structured disentangled representations; one possible direction is to encourage group independence. We leave this for future work.

## 5 Acknowledgement

Dejiao Zhang and Laura Balzano were supported by DARPA grant 16-43-D3M-FP-037. The authors also would like to thank Zeyu Sun for the initial participation in this project.

<sup>3</sup>Although the truth discrete factor is provided, we evaluate the disentanglement quality only in terms of the continuous representations since the pixel-wise difference between different categories are very small.

## References

- Thomas M Cover and Joy A Thomas. *Elements of information theory*. John Wiley & Sons, 2012.
- Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- Eric Jang, Shixiang Gu, and Ben Poole. Categorical reparameterization with gumbel-softmax. *arXiv preprint arXiv:1611.01144*, 2016.
- Satosi Watanabe. Information theoretical analysis of multivariate correlation. *IBM Journal of research and development*, 4(1):66–82, 1960.
- Tian Qi Chen, Xuechen Li, Roger Grosse, and David Duvenaud. Isolating sources of disentanglement in variational autoencoders. *arXiv preprint arXiv:1802.04942*, 2018.
- Babak Esmaeili, Hao Wu, Sarthak Jain, Siddharth Narayanaswamy, Brooks Paige, and Jan-Willem van de Meent. Hierarchical disentangled representations. *arXiv preprint arXiv:1804.02086*, 2018.
- Irina Higgins, Loic Matthey, Arka Pal, Christopher Burgess, Xavier Glorot, Matthew Botvinick, Shakir Mohamed, and Alexander Lerchner. beta-vae: Learning basic visual concepts with a constrained variational framework. 2016.
- Shengjia Zhao, Jiaming Song, and Stefano Ermon. Infovae: Information maximizing variational autoencoders. *arXiv preprint arXiv:1706.02262*, 2017.
- Emilien Dupont. Joint-vae: Learning disentangled joint continuous and discrete representations. *arXiv preprint arXiv:1804.00104*, 2018.
- Felix Vsevolodovich Agakov. *Variational Information Maximization in Stochastic Environments*. PhD thesis, University of Edinburgh, 2005.
- Takeru Miyato, Andrew M Dai, and Ian Goodfellow. Virtual adversarial training for semi-supervised text classification. 2016.

## A Supplement for the numerical results

$$\begin{aligned}
 \mathcal{L}_{\text{VAE}} &= \mathbb{E}_{p(\mathbf{y}, \mathbf{z}|\mathbf{x})} [q(\mathbf{x}|\mathbf{y}, \mathbf{z})] - D_{\text{KL}}(p(\mathbf{z}|\mathbf{x})||r(\mathbf{z})) - D_{\text{KL}}(p(\mathbf{y}|\mathbf{x})||r(\mathbf{y})) \leftarrow \text{ELBO} \\
 &= \underbrace{\mathbb{E}_{p(\mathbf{y}, \mathbf{z}|\mathbf{x})} [q(\mathbf{x}|\mathbf{y}, \mathbf{z})]}_{\textcircled{1}} - \underbrace{I(\mathbf{x}; \mathbf{y})}_{\textcircled{2}} - \underbrace{D_{\text{KL}}(p(\mathbf{y})||r(\mathbf{y}))}_{\textcircled{3}} - \underbrace{I(\mathbf{x}; \mathbf{z})}_{\textcircled{4}} - \underbrace{D_{\text{KL}}(p(\mathbf{z})||r(\mathbf{z}))}_{\textcircled{5}} \\
 \beta\text{-VAE: } &\textcircled{1} - \beta (\textcircled{2} + \textcircled{3}) - \beta (\textcircled{4} + \textcircled{5}) \quad \text{InfoVAE: } \textcircled{1} - \beta \textcircled{3} - \beta \textcircled{5} \\
 \text{Joint-VAE: } &\textcircled{1} - \beta |\textcircled{2} + \textcircled{3} - C_{\mathbf{y}}| - \beta |\textcircled{4} + \textcircled{5} - C_{\mathbf{z}}|
 \end{aligned}$$

Figure 4: Summarization of relevant works.  $\beta$ -VAE modifies ELBO by increasing the penalty on the KL divergence terms. Both InfoVAE and WAE drop the mutual information terms from ELBO. JointVAE seeks to control the mutual information by pushing their upper bounds, *i.e.*, the associated KL divergence terms, towards progressively increased target values,  $C_{\mathbf{y}}$  and  $C_{\mathbf{z}}$ .

Table 1: Encoder and Decoder architecture for MNIST and Fashion MNIST.

Encoder	Decoder
Input vectorized $28 \times 28$ grayscale image	Input $\mathbf{y} \in \mathbb{R}^{10}$ and $\mathbf{z} \in \mathbb{R}^{10}$
FC. 500 BatchNorm ReLU	FC. 500 ReLU
FC. $2 \times 500$ BatchNorm ReLU	FC. 500 ReLU
FC. $20 (\mu_{\mathbf{z}}, \log \sigma_{\mathbf{z}}) + 10 (p_{\mathbf{y}})$	FC. $28 \times 28$ Sigmoid

Table 2: Encoder and Decoder architecture for dSprites.

Encoder	Decoder
Input vectorized $64 \times 64$ grayscale image	Input $\mathbf{y} \in \mathbb{R}^3$ and $\mathbf{z} \in \mathbb{R}^{10}$
FC. 1200 ReLU	FC. 1200 ReLU
FC. 1200 ReLU	FC. 1200 ReLU
FC. $2 \times 1200$ ReLU	FC. 1200 ReLU
FC. $20 (\mu_{\mathbf{z}}, \log \sigma_{\mathbf{z}}) + 3 (p_{\mathbf{y}})$	FC. $28 \times 28$ Sigmoid