

# Conditional Variational Autoencoders with Fuzzy Inference

Yury Gurov<sup>1</sup>[0000-0002-7033-9996] and Danil Khilkov<sup>1</sup>[0000-0001-9284-6924]

NIIAS Institute of Informatization, Automation and Communication in Railway Transport, Russia, Moscow 109029 Nizhegorodskaya str., 27 bldg. 1 [info@vniias.ru](mailto:info@vniias.ru)  
[www.vniias.ru/](http://www.vniias.ru/)

**Abstract.** We present an approach to constructing Conditional Variational Autoencoders (C-VAE) models with fuzzy inference during classification. This preserves disentangling capabilities of VAE and at the same time performs latent space clusterization. Fuzzy C-VAE model provides useful features for anomaly detection, utilizing partially labeled datasets and controlled generation of new samples.

**Keywords:** fuzzy logic · deep learning · fuzzy inference · Conditional Variational Autoencoders · fuzzy cvae · neuro-fuzzy

## 1 Introduction

Hybrid neuro-fuzzy systems has a long time history and is still an active research area [?]. Main feature that attract attention to neuro-fuzzy systems is possibility to combine the power of neural network with advantages of fuzzy logic, such a human-like reasoning. At this work we propose an approach to constructing Conditional Variational Autoencoders (C-VAE) [?, ?, ?] models with fuzzy inference during classification phase. This approach may preserve disentangling capabilities of VAE and at the same time performs latent space clusterization. Such fuzzy C-VAE model provides useful features for anomaly detection, utilization of partially labeled datasets and controlled generation of new samples.

The source code is available at GitHub repository (<https://github.com/kenoma/pytorch-fuzzy>).

## 2 Related work

In [?] attempt to apply fuzzy logic to the latent space of VAE was made with fuzzy c-mean clustering. Main drawback of this approach is that it requires prior knowledge about problem domain number of clusters and their interpretations.

In [?] conditional VAE was modified in a way to process partially observed datasets. Authors proposed method that augments the conditional VAEs with a prior distribution for the missing covariates and estimates their posterior using amortised variational inference. At first sight this approach has nothing to do with fuzzy logic, but it provides insight into the problem of latent space clustering.

### 3 Methods

#### 3.1 Variational Autoencoders

Variational inference is used to approximate a posterior distribution of a directed graphical model whose latent variables and parameters are intractable. The Variational Auto-Encoder (VAE) combines this approach with an autoencoder framework to learn the prior distribution of a latent space,  $p_\theta(z)$ , with parameters  $\theta$ . The idea is that the prior distribution can then be sampled to produce a latent code,  $z$ , which is passed as input to the decoder to produce a sample output,  $\tilde{x}$ . VAEs consists of two NN for the probabilistic encoding and decoding process (see Figure 1a). As the true underlying distribution of the posterior is intractable and complex, a simple parametric surrogate distribution,  $q_\Phi(z|x)$  (such as a Gaussian), with parameters  $\Phi$ , is assumed to approximate the distribution and is optimized for best fit. The encoder network implicitly models the surrogate distribution, by mapping the distribution parameters,  $\Phi$ , during the training process. The resulting model,  $q_\Phi(z|x)$ , is referred to as the recognition model. The optimization process of the recognition model revolves around minimizing the Kullback-Leibler (KL) divergence between the posterior and surrogate distributions. Once the latent prior distribution is learned,  $z$  can be sampled via the reparameterization trick. The (probabilistic) decoder network performs a mapping of the latent code to a structured sample output for each sample, thus producing a distribution of outputs,  $p_\theta(x|z)$ .

#### 3.2 Conditional VAE

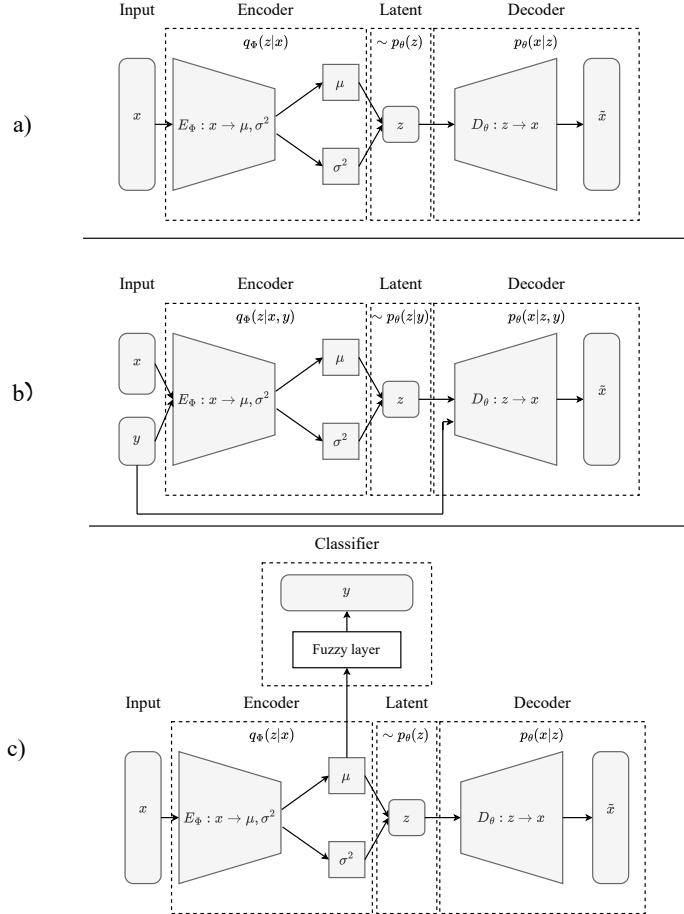
The C-VAE expands upon the framework of the VAE, by combining variational inference with a conditional directed graphical model. In the case of C-VAE, the objective is to learn a prior distribution of the latent space that is conditioned on an input variable  $y$  such that  $p_\theta(z|y)$  (see Fig. ??b). The conditioning of the distributions results in a prior that is modulated, by the input variable, creating a method to control modality of the output.

#### 3.3 Fuzzy C-VAE

We propose C-VAE architecture where additional conditions are applied only to  $/mu$  component in order to reorganize the latent space structure (see Fig. ??c). Reorganization achieved by using fuzzy term functions, where each term associated with sole condition i.e. label. Multidimensional Gaussian function is used to represent the fuzzy term function:

$$\nu(z, A_i) = e^{\frac{1}{2} \|\tilde{z}\|_i^2},$$

where  $m$  is a latent space dimension size,  $i$  denotes term number,  $\tilde{z} = [z_1, z_2, \dots, z_m, 1]$  and  $A_i$  is transformation matrix in form



**Fig. 1.** Overview of the a) VAE b) CVAE by [?] and c) proposed fuzzy C-VAE model.

$$A_{(m+1) \times (m+1)} = \begin{bmatrix} s_1 & a_{12} & \cdots & a_{1m} & c_1 \\ a_{21} & s_2 & \cdots & a_{2m} & c_2 \\ \vdots & \vdots & \ddots & \vdots & c_3 \\ a_{m1} & a_{m2} & \cdots & s_m & c_m \\ 0 & 0 & \cdots & 0 & 1 \end{bmatrix},$$

with  $c_{1\dots m}$  centroid position,  $s_{1\dots m}$  scaling factor and  $a_{1\dots m, 1\dots m}$  as alignment coefficients. Such representation of multidimensional Gaussian term function easily can be adopted for use in any modern machine learning framework. Set of  $\nu(z, A_i)$  we call a fuzzy layer. Intuition behind fuzzy layer is that during training procedure every input vector  $z$  will be forced to group closer near centroid

of corresponding term function. Disentangling features of VAE combined with clustering possibilities of fuzzy layer provides a way to learn supervised latent space which can be useful for anomaly detection and other tasks we discuss further.

### 3.4 Learning Fuzzy C-VAE

To train fuzzy C-VAE we use the same loss function as in standard VAE with addition of fuzzy layer loss:

$$\text{Loss} = \text{MSE}(\tilde{x}, x) + \text{KL}(\mu, \log \sigma^2) + \text{FZ}(\tilde{y}, y),$$

where  $\text{MSE}(\tilde{x}, x)$  is reconstruction loss,  $\text{KL}(\mu, \log \sigma^2)$  is the KL-divergence (for more details see [?]) and  $\text{FZ}(\tilde{y}, y)$  represents the mean squared error between the output and target conditional vector.

Main drawback of fuzzy layer is that in high dimensional cases it's hard to find good initial values for centroids and scaling factors mainly due to vanishing gradients. In such a case it is possible to pass to fuzzy layer subsection of vector  $/mu$  leaving remained part to be trained by VAE without any conditional restrictions.

## 4 Experiments

In this paper we would like to demonstrate ideas of fuzzy C-VAE on playground MNIST dataset. To make demonstration fancy we provide additional label to samples of MNIST dataset. This label separates numbers with closed round loops in outline 0, 6, 8, 9 from numbers without it 1, 2, 3, 4, 5, 7. To make reasonable reconstruction loss we set size of the latent vector equal to 12 but only 2 first values of this vector are passed to fuzzy layer. The fuzzy CVAE model is trained using the Adam optimizer [?].

### 4.1 Latent space clusterization

On Figure ?? is depicted latent space structure resulted by vanilla VAE without any conditional restrictions. Despite the fact that points in VAE latent space are grouped in clusters corresponding to each number this structure has very complex topology. Without application of prior knowledge about number labels task of extracting corresponding clusters is very challenging.

Passing label information directly to fuzzy C-VAE during training leads to fine grained latent space structure as shown on Figure ???. First two components of latent vector on which fuzzy layer is applied have formed clusters corresponding to each number. Meanwhile remained components of latent vector preserves complex topology like pure VAE. Reconstruction losses for VAE and fuzzy C-VAE during our experiments were almost the same while KL-loss for fuzzy C-VAE was slightly higher all the time. Classification accuracy of fuzzy

C-VAE depends on many factors of network topology and training scenario but limit of 99% is easily achieved. Figure ?? shows how fuzzy C-VAE is able to classify samples.

#### 4.2 Controlled samples generation

After training procedure it is possible extract from matrices  $A_i$  cluster characteristics such a centroids, scaling and alignment factors to understand latent space structure. Thats make possible use fuzzy C-VAE for sample generation with predefined characteristics as shown on Figure ??.

#### 4.3 Anomaly detection

Fuzzy C-VAE provides convinient way to design anomaly detection model. First of all every term at fuzzy layer defines a subdomain in latent space which is described through multidimentional Gaussian.

#### 4.4 Learning on partially labeled dataset

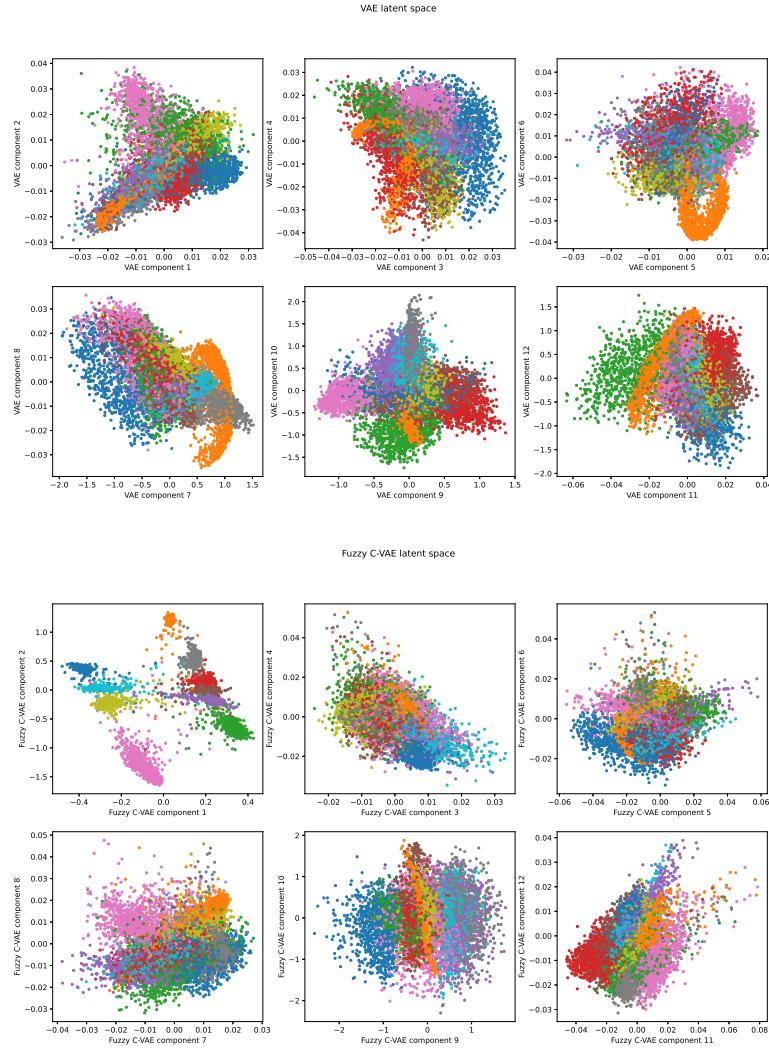
### 5 Discussion

In this paper, we introduced a novel fuzzy inference layer to improve the performance of conditional VAEs. We achieve this by making trainable multidimensional representation of fuzzy term. The method that we proposed is applicable to a variety of conditional VAE models. The efficacy of our proposed method was demonstrated on MNIST dataset. Whereas fuzzy conditions influence on VAE should be discussed more deeply we leave it for future work.

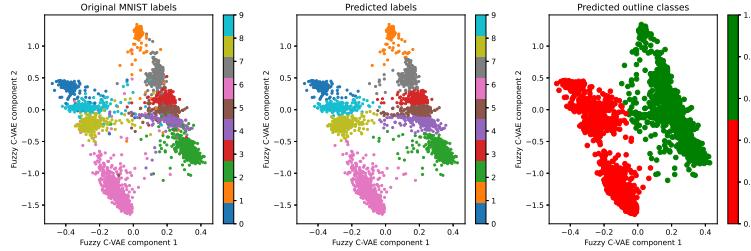
### References

1. Bölat, K., Kumbasar, T.: Interpreting variational autoencoders with fuzzy logic: A step towards interpretable deep learning based fuzzy classifiers. In: 2020 IEEE International Conference on Fuzzy Systems (FUZZ-IEEE). pp. 1–7 (July 2020). <https://doi.org/10.1109/FUZZ48607.2020.9177631>
2. de Campos Souza, P.V.: Fuzzy neural networks and neuro-fuzzy networks: A review the main techniques and applications used in the literature. *Applied Soft Computing* **92**, 106275 (2020). <https://doi.org/10.1016/j.asoc.2020.106275>, <https://www.sciencedirect.com/science/article/pii/S1568494620302155>
3. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization (2017). <https://doi.org/10.48550/arXiv.1412.6980>
4. Kingma, D.P., Welling, M.: An introduction to variational autoencoders. *Foundations and Trends® in Machine Learning* **12**(4), 307–392 (2019). <https://doi.org/10.1561/2200000056>, <http://dx.doi.org/10.1561/2200000056>
5. Kingma, D.P., Welling, M.: Auto-encoding variational bayes (2022). <https://doi.org/https://doi.org/10.48550/arXiv.1312.6114>

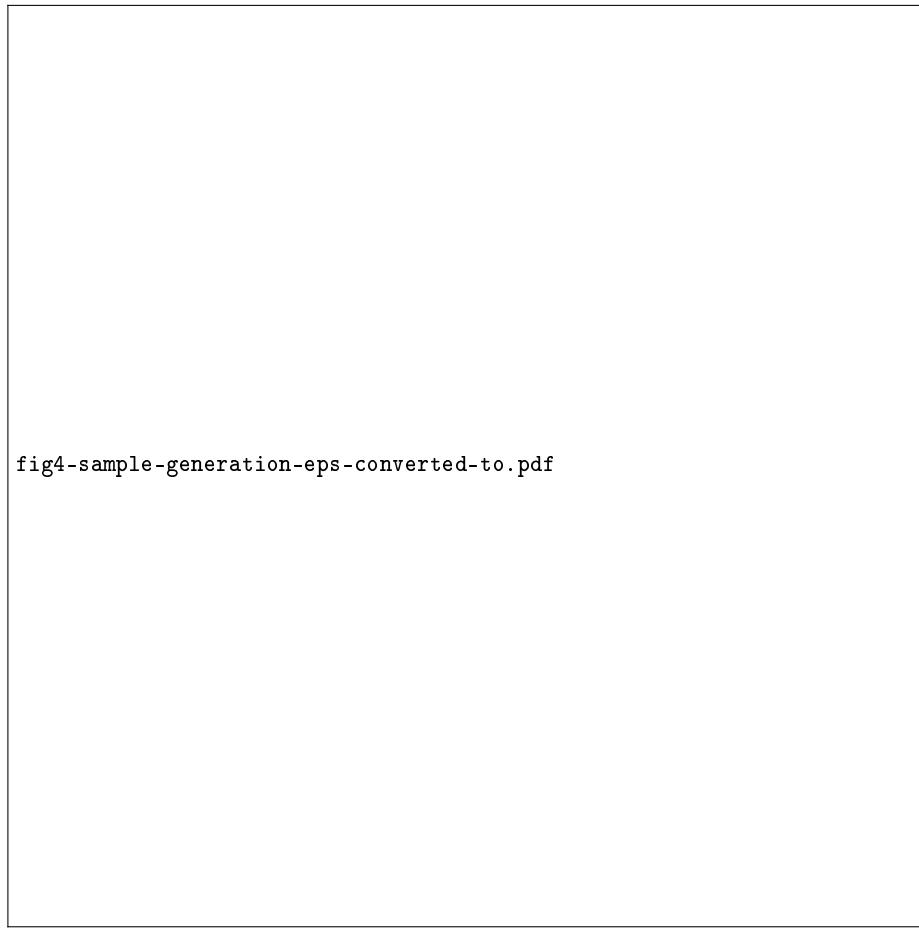
6. Ramchandran, S., Tikhonov, G., Lönnroth, O., Tiikkainen, P., Lähdesmäki, H.: Learning conditional variational autoencoders with missing covariates. *Pattern Recognition* **147**, 110113 (2024). <https://doi.org/https://doi.org/10.1016/j.patcog.2023.110113>, <https://www.sciencedirect.com/science/article/pii/S0031320323008105>
7. Sohn, K., Yan, X., Lee, H.: Learning structured output representation using deep conditional generative models. In: Proceedings of the 28th International Conference on Neural Information Processing Systems - Volume 2. p. 3483–3491. NIPS’15, MIT Press, Cambridge, MA, USA (2015). <https://doi.org/10.5555/2969442.2969628>



**Fig. 2.** Latent space granulation for vanila VAE (top) and Fuzzy C-VAE (bottom)



**Fig. 3.** Fuzzy C-VAE latent space colored by true number labels (left) predicted number labels (center) and predicted outline class (right)



**Fig. 4.** Fuzzy C-VAE latent space colored by true number labels (left) predicted number labels (center) and predicted outline class (right)