

---

# EXAM PROJECT FOR PML 2022/2023

---

## REPORT

**Christian Dybdahl Troelsen**  
Department of Computer Science  
University of Copenhagen  
Universitetsparken 1  
DK-2100 Copenhagen Ø  
tfp233@alumni.ku.dk

**Jens Sørensen**  
Department of Computer Science  
University of Copenhagen  
Universitetsparken 1  
DK-2100 Copenhagen Ø  
qrw992@alumni.ku.dk

**Mathias Rasmussen**  
Department of Computer Science  
University of Copenhagen  
Universitetsparken 1  
DK-2100 Copenhagen Ø  
tjc725@alumni.ku.dk

January 18, 2023

## 1 Density modeling

### 1.1 Implement a convolutional VAE

**Architecture of the convolutional VAE** Our implementation of the convolutional VAE is similar in architecture to the original VAE. The primary difference between the two models is that the two linear layers in the encoder and decoder part of the original VAE have been replaced by respectively two convolutional layers and two transposed convolutional layers in the convolutional VAE. The two convolutional layers in the encoder part of the convolutional VAE both use kernels of size  $3 \times 3$ ,  $1 \times 1$  padding and a stride of  $2 \times 2$ . The first convolutional layer has 1 input channel and 16 output channels, while the second convolutional layer has 16 input channels and 32 output channels. The transposed convolutional layers in the decoder part of the convolutional VAE complement the aforementioned convolutional layers. Both transposed convolutional layers use kernels of size  $3 \times 3$ ,  $1 \times 1$  padding  $1 \times 1$ , output padding and a stride of  $2 \times 2$ . The first transposed convolutional layer has 32 input channels and 16 output channels, while the second transposed convolutional layer has 16 input channels and 1 output channel. Conceptually, the encoder part of the convolutional VAE condenses its input images  $\mathbf{x}$  of size  $D = N \times N$  into a compact but feature rich representation  $\mathbf{z}$  of size  $C = 2$ , while the corresponding decoder "unfolds" this compact representation in reverse order, ultimately producing output images  $\mathbf{y}$  of the same size as  $\mathbf{x}$ .

**Parameter estimation** In order to optimize the convolutional VAE model we estimate the parameters  $\phi$  and  $\theta$  that maximize the ELBO, which is a lower bound on the log likelihood of the data  $\mathbf{x}$  and is defined as  $\mathcal{L}(\theta, \phi, \mathbf{x}) = \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\ln p_\theta(\mathbf{x}|\mathbf{z})] - \text{KL}(q_\phi(\mathbf{z}|\mathbf{x})||p_\theta(\mathbf{z}))$ . This is the same criterion used to optimize the original VAE. In particular, both models use the reparameterization trick to sample latent variables  $\mathbf{z}$  from a factorized gaussian approximate posterior  $q_\phi(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mathbf{z}; \mathbf{\mu}, \text{diag}(\sigma^2))$ , so the Kullback-Leibler divergence, which acts as a regularization term in the ELBO, simplifies to  $\frac{1}{2} \sum_{d=1}^D (\sigma_d^2 + \mu_d^2 - 1 - \ln(\sigma_d^2))$ , assuming that the prior over  $\mathbf{z}$  is also a diagonal gaussian  $p_\theta(\mathbf{z}) = \mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbf{I})$ . Here  $\ln \sigma^2$  and  $\mu$  are the outputs of the encoder when applied to  $\mathbf{x}$ . In order to compute the first term in the ELBO, which is the expected likelihood of  $\mathbf{x}$  given  $\mathbf{z}$ , we utilize the output  $\mathbf{y}$  of the decoder. We feed  $\mathbf{y}$  through a sigmoid activation layer, which produces the mean parameters  $\mathbf{p}$  of a factorized multivariate Bernoulli distribution  $\prod_j^D \text{Bernoulli}(x_j; p_j)$  acting as the likelihood  $p_\theta(\mathbf{x}|\mathbf{z})$  of the input  $\mathbf{x}$  given latent variable  $\mathbf{z}$ . We then approximate the expected log likelihood by computing the negative of the binary cross entropy of  $\mathbf{p}$  and  $\mathbf{x}$ , where  $\mathbf{x}$  is treated as a set of probability values. Given the simplified ELBO, we optimize the parameters of both models using a variant of gradient descent, namely the Adam algorithm with a learning rate of 0.001 and a batch size of 128. In order to ensure convergence we run the Adam algorithm for a total of 25 epochs for each model. We train both models on the complete MNIST training set.

**Performance of VAE models** We compare the performance of the two VAE models on the MNIST test set using two quantitative measures and three qualitative experiments. The two quantitative measures are the mean log likelihood and the MSE loss. The three qualitative experiments are

1. Clustering MNIST test data in latent space by mapping each observation  $x_i$  to corresponding mean parameters  $\mu_i$  using the encoder and subsequently plotting these  $\mu_i$  with colors based on the actual labels  $t_i$  associated with each  $x_i$ .
2. Exploring the latent space by sampling latent variables  $z$  deterministically from a regular grid, then using the decoder to map these  $z$  to corresponding mean parameters  $p$ , and finally visualizing these  $p$  as images.
3. Reconstructing selected observations  $x_i$  from the MNIST test set by mapping these to posterior distributions  $q_\phi(z_i|x_i)$ , then randomly sampling  $z_i$  from these posterior distributions, and finally using the decoder to map these  $z_i$  to corresponding mean parameters  $p_i$ , which are finally interpreted as reconstruction of the inputs  $x_i$ .

A comparison of the mean log likelihood and MSE loss for the two VAE models is shown in Table 1. As can be seen,

## 1.2 Alternative models

### 1.2.1 Probabilistic PCA

### 1.2.2 Gaussian Mixture Model

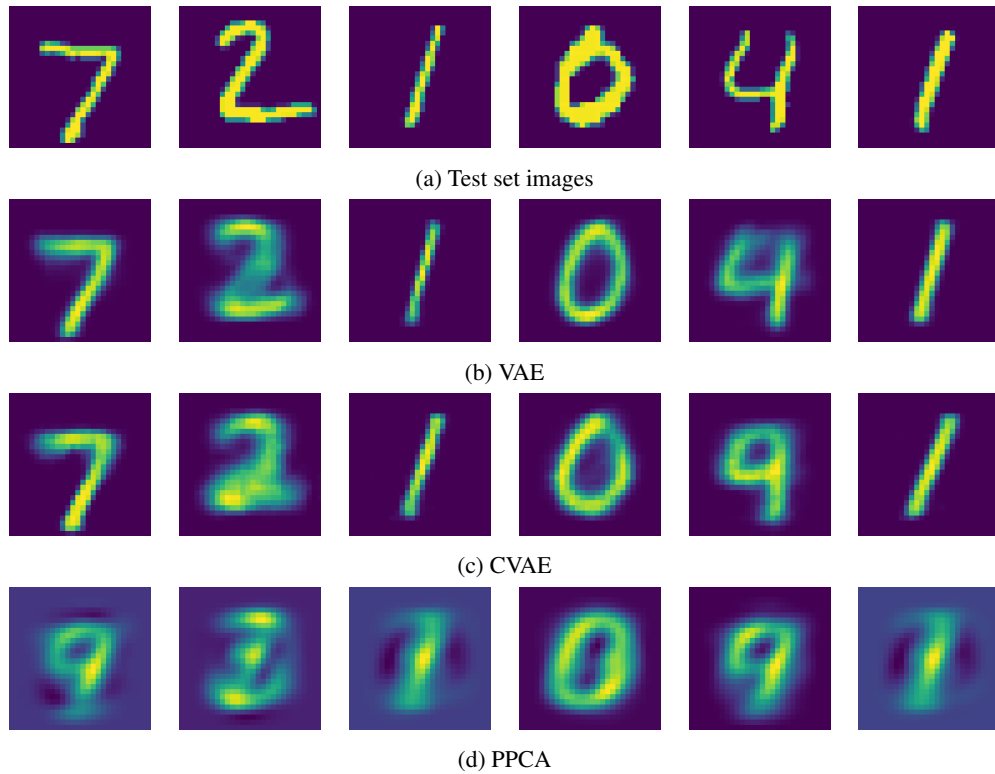


Figure 1: Comparison of MNIST test set images and corresponding mean parameters generated by density models

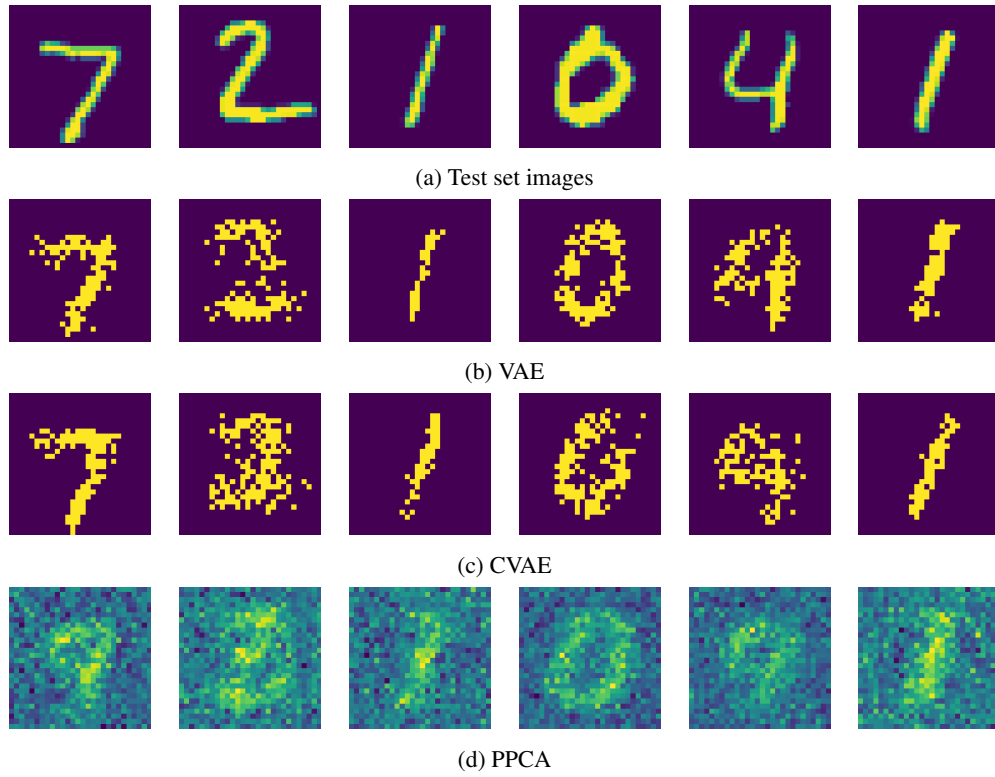
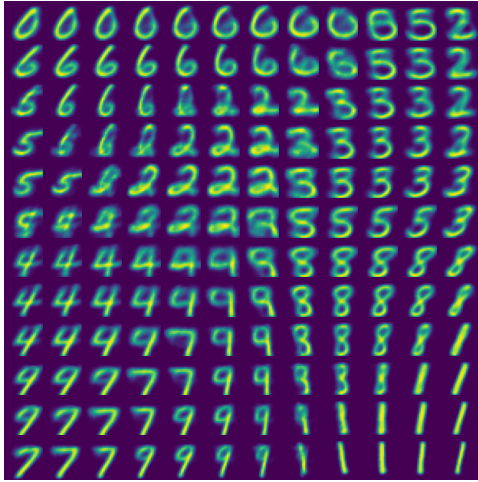
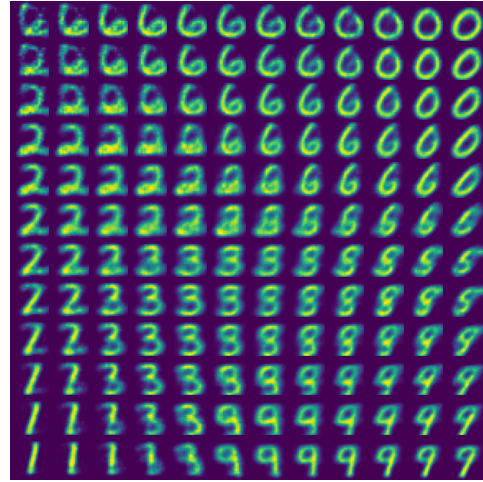


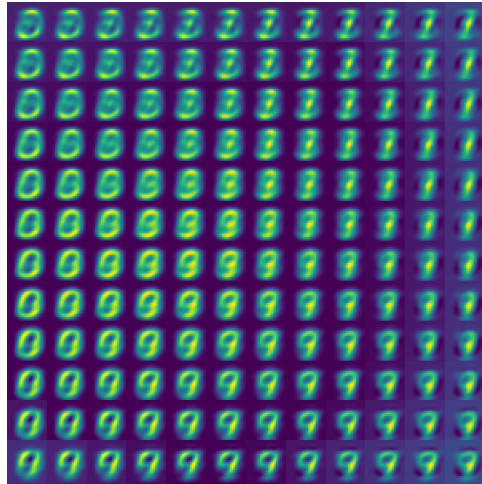
Figure 2: Comparison of MNIST test set images and corresponding reconstructions sampled from density models



(a) VAE



(b) CVAE



(c) PPCA

Figure 3: Interpolating images from latent space variables using trained density models.

	Log-Likelihood/ELBO	MSE
<b>VAE</b>	-1.428134e+02	0.037524
<b>CVAE</b>	-1.568176e+02	0.044254
<b>PPCA</b>	-4.329656e+03	3629.250732
<b>GMM</b>	-1.067011e+07	3831.718461

Table 1: Model performance metrics

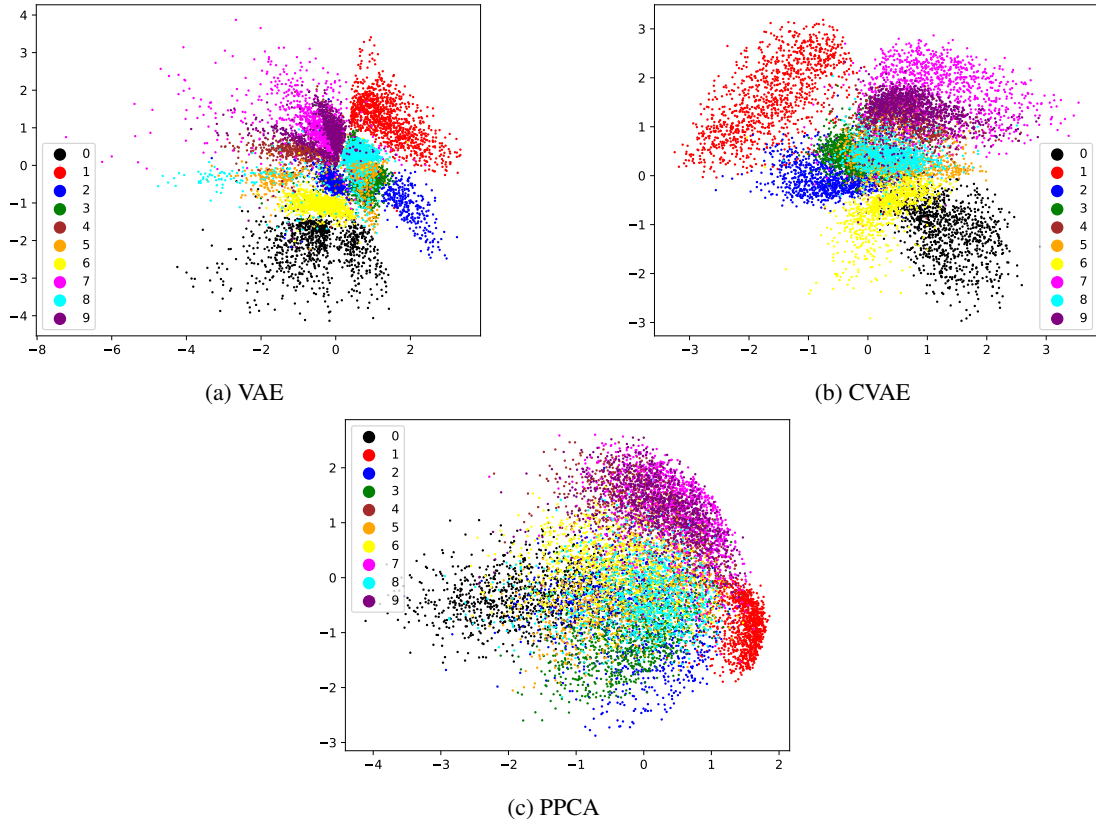


Figure 4: Clustering on MNIST test (projection to latent space) using trained density models.

## 2 Function fitting

### 2.1 Fitting a GP with Pyro

### 2.2 Bayesian Optimization

## 3 Bibliography