

# AIML 425 Assignment 3

James Thompson  
Victoria University of Wellington  
300680096

## I. INTRODUCTION

Autoencoders are a type of neural network used for unsupervised learning of efficient latent space representations of data. An auto encoder consists of two main parts the encoder which compresses the input data into a lower dimensional space and the decoder which reconstructs the original data. Variational Autoencoder are improved Autoencoders that change the output of the encoder to be a distribution over the latent space rather than a fixed vector and uses a ELBO loss function. I will explore both the generative capabilities of regular Autoencoders and Variational Autoencoders as well as the meaning of the latent space and how much information is retained in the latent space. I will train both models on a generated dataset of shapes of different size and locations on a 28x28 canvas. Both models achieve good reconstruction of the input data, however only the VAE is able to generate reasonable new images by sampling the latent space.

## II. THEORY

Autoencoder and their variations are a type of neural networks architecture that are used for unsupervised learning of an efficient representation of the input data. Basic autoencoders use a encoder-decoder architecture to compress and then reconstruct the input data. Variational Autoencoders (VAEs) extend this architecture by introducing a probabilistic approach to the encoding process, allowing for the latent space to be a known distribution. This section provides a description of both model types as well as a discussion of the latent space.

### A. Autoencoders

TODO: add mathematical/probabilistic description of model

An Auto encoder has three important parts the encoder  $f_{enc}$ , the decoder  $f_{dec}$  and the latent layer  $z$ . The latent layer is typically of a much lower dimension (2D-32D) than the input data. The encoder and decoder are typically implemented as neural networks with learnable parameters  $\phi$  and  $\theta$ . Furthermore the encoder and decoder are generally the inverse of each other.

$$z = f_{enc}(x; \phi) \quad (1)$$

$$\hat{x} = f_{dec}(z; \theta) \quad (2)$$

The input data  $x$  goes through  $f_{enc}$  to  $z$  then through  $f_{dec}$  to produce the reconstruction  $\hat{x}$ . This model is then trained based on the reconstruction loss between the input and output data.

$$\mathcal{L}_{AE} = \|x - \hat{x}\|^2 \quad (3)$$

$$(4)$$

TODO: update the loss to include mention of the variance control / MMD.

Training with this loss function means that the model will learn a compressed representation of the data. However there is no reason that the distribution of the latent space will be smooth or continuous (i.e that all squares will be near each other in the latent space). This means that the generative capabilities of the model will be limited. Only certain part of the latent space will correspond to valid images, without access to the encoder and training data we have no way of knowing which parts of the latent space are valid. This means that sampling the latent space to generate new images is likely to produce random noise rather than valid images.

### B. Variational Autoencoders

Variational Autoencoders (VAEs) are a type of generative model that builds upon the architecture of traditional autoencoders by introducing a probabilistic approach to the encoding process. In a VAE, the encoder does not produce a single point in the latent space but rather a distribution over the latent space, typically modeled as a Gaussian distribution. The goal is to learn the parameters of the generative model  $p_{\theta}(x|z)$  which would allow us to generate new data points by sampling from the latent space.

$$q(z|x) = \mathcal{N}(z; \mu(x), \text{diag}(\sigma^2(x))) \quad (5)$$

$$\mu(x), \sigma(x) = f_{enc}(x; \phi) \text{ where } \mu \text{ is mean and } \sigma \text{ is variance}$$

$$z \sim q(z|x) \quad (6)$$

$$\hat{x} = f_{dec}(z; \theta) \quad (7)$$

This changes the latent space so that it comes from a normal distribution. Therefore the latent space is continuous, smooth and has a known distribution. The latent space will now be valid everywhere in the space, meaning that we can sample from the latent space to generate new images.

Furthermore, VAEs are trained using a loss function that combines the reconstruction loss (similar to traditional autoencoders) with a regularization term that encourages the learned latent distributions to be close to a prior distribution (typically a standard normal distribution). This regularization is achieved using the Kullback-Leibler (KL) divergence, which measures the difference between the learned distribution and the prior distribution. The overall loss function for a VAE can be expressed as:

### III. EXPERIMENTS

$$\mathcal{L}_{VAE} = \mathcal{L}_{reconstruction} + \beta \mathcal{L}_{KL} \quad (8)$$

Using MSE for reconstruction loss and Gaussian prior:

$$\begin{aligned} &= \mathbb{E}_{x \sim p_{data}(x)} [\log p_{\theta}(x|z)] - \beta D(q(z|x)||p(z)) \\ &= \|x - \hat{x}\|^2 + \beta \frac{1}{2} \sum_{i=1}^d (\mu_i^2 + \sigma_i^2 - \log(\sigma_i^2) - 1) \end{aligned} \quad (9)$$

This loss function is derived using the Jensen's inequality to maximize the evidence lower bound (ELBO) of the data likelihood, which will also maximize the likelihood of  $\theta$ . The reconstruction loss ensures that the decoded outputs are similar to the inputs, while the KL divergence term regularizes the latent space to follow the prior distribution. The hyperparameter  $\beta$  controls the trade-off between these two objectives.

For back propagation through the stochastic sampling process, VAEs employ the reparameterization trick. This involves expressing the latent variable  $z$  as a deterministic function of the encoder outputs and a random noise variable  $\epsilon$  sampled from a standard normal distribution:

$$z = \mu + \sigma \cdot \epsilon, \quad \epsilon \sim \mathcal{N}(0, I) \quad (10)$$

Where  $\epsilon$  is sampled for the minibatch and assumed to be constant for the back-propagation step. This allows gradients to be computed with respect to the encoder parameters, enabling effective training of the VAE using standard back-propagation techniques.

#### C. Latent space information

1) *Information rate*: When we compress data, we inevitably lose some information. The amount of information retained in the compressed representation can be quantified using the concept of information rate. Assuming that we add noise to the latent space,  $z = y + \epsilon$  where  $\epsilon \sim \mathcal{N}(0, \sigma^2)$ , the information rate can be calculated as:

$$R = \sum_i \frac{1}{2} \log_2 \left( 1 + \frac{\sigma_{i,Z}^2}{\sigma_{i,\epsilon}^2} \right) \quad (11)$$

Where  $\sigma_{i,Z}^2$  is the variance of the  $i$ -th dimension of the latent space and  $\sigma_{i,\epsilon}^2$  is the variance of the noise added to that dimension.

This means that for a given noise level, we can estimate the amount of information that can be retained in the latent space by sampling the encoder outputs and measuring the variance of each dimension.

2) *Information meaning*: The information that the latent space retain can be interpreted in multiple ways. One way is that each dimension of the latent space could represent a specific feature of the input data. For example in my data generation process I generate each shape with only 4 numbers (x position, y position, size, shape). Therefore one could imagine a 4D latent space that perfectly encodes these 4 features. In practice this is unlikely to happen, but we can still interpret the latent space as a compressed representation of the input data.

#### A. Dataset

The dataset used for training both models is a generated dataset of simple shapes (circles, squares, triangles) on a 28x28 canvas. Each shape is generated with random size and position. The shape may go up to the bounds but will never go out it and be clipped. There are 8 sizes ([7, 14]) which means that there are a total 8002 unique images. The bounds of sizes were determined by how small the circles still look circle (lower bound) and no taking up more than half of the canvas (upper bound). The dataset used for training is 90% of the total possible images, and the test set is the remaining 10%. Holding a validation allows me to test the models ability to reconstruct unseen images, and so not just memorize the training data.

The size of the square determines its width/height, the size of the circle determines its diameter and the size of the triangle determines the height (but the size is modified to be a bit shorter to make it fit with the other image sizes). The shapes are drawn in white on a black background. The triangle is an upright isosceles triangle.

#### B. Variational autoencoder

1) *setup*: The VAE model is two fully connected neural networks. The encoder outputs 2 values (mean and variance) while the decoder outputs a single scalar. The loss function used was the ELBO loss function with a  $\beta$  of 0.65 (using BCE for reconstruction error and KL divergence for regularization). A complete list of hyperparameters is shown in table II. The fine tuning of parameters was done through trial and error using the loss curve and generative performance metrics as a guide.

2) *generative performance*: There are different ways of understanding the generative performance of the model. This situation we are involves a finite set of valid images (8002) in a very large space of images ( $2^{784}$ ). Using this valid image set we can calculate three different metrics. Therefore we can use this valid image set to evaluate; how many of the valid images can the model generate (coverage), how close are the generated images to the nearest valid image (nearest neighbor performance) and how similar is the distribution of generated images to the distribution of valid images (KL divergence).

$$\text{coverage} = \frac{|\{x \in \text{images} \mid \min_{\hat{x} \in \text{generated}} d(x, \hat{x}) < \text{thres}\}|}{|\text{images}|}$$

$$\text{nearest neighbor} = \frac{1}{|\text{images}|} \sum_{x_i \in \text{images}} \min_{\hat{x}_j \in \text{generated}} d(x_i, \hat{x}_j)$$

$$\begin{aligned} \text{KLD} &= D_{KL}(p_{data} || p_{model}) \\ &\approx \sum_{x_i \in \text{images}} p_{data}(x_i) \log \frac{p_{data}(x_i)}{p_{model}(x_i)} \\ &= \sum_{x_i \in \text{images}} \frac{1}{|\text{images}|} \log \frac{1}{q_{counts}(x_i)} \end{aligned}$$

Where  $d(x_i, \hat{x}_j)$  is the distance between two images, and the thres is a hyperparameter that determines how close the

generated image needs to be to a valid image to count as "an attempt at generating it", and  $q_{counts}(x_i)$  is the number of times the model has generated an image close enough to  $x_i$  (within the threshold). I choose euclidean as the distance metric and the threshold of 3 by looking at sample images and aligning the cutoff and distance metric with my intuition.

The results of the generative performance can be found in I

### C. Autoencoder

1) *setup*: The Autoencoder model is a simpler version of the VAE model. The encoder and decoder are both just feedforward networks but unlike the VAE the encoder outputs a single vector rather than a distribution. The loss function used was just the mean squared error between the input and reconstruction input. Additionally a simple variance control was added that pushes the latent space to have unit variance. Lastly some gaussian noise was randomly added to the latent space. A complete list of hyperparameters is shown in table II. Specifically compared to the VAE, the model was trained for longer, is twice the size and uses dropout. The fine tuning was done through trial and error using the loss curve and generative performance. Of key importance was pushing the network so that it could actually generate anything resembling a valid image. To do this required both the noise in the latent space and dropout in the encoder.

2) *Information rate*: The information rate of the latent space was calculated by getting the encoder outputs from all of the possible images. Then using equation 11 to calculate the information rate. The estimated information was 12.9 bits. This means that the latent space could represent 7604 distinct images. This is quite close to the 8002 valid images and would fit into a half float. Which achieves a good compression of the data into the equivalent space of less than 2 bytes.

3) *Information meaning*: There are two methods that I will use to understand the latent space. The first is to check for correlation between known features (x position, y position, size, shape) and the latent dimensions. The second is to change the latent dimensions and see how the reconstruction changes. This can be done by either changing one dimension at a time or zeroing out one dimension at a time.

The correlation between known features and latent dimensions show 2 that there is some correlation between some latent dimensions. Particularly the size both have some correlations around 0.3, other latent have some trivial correlations 0.1 but nothing strong. Furthermore looking at the distribution of latent dimensions grouped by shape 1 shows that there is no strong correlation between shape and any latent dimension. This means that the model has not learned to dedicate a specific dimension to shape. This is more or less what we would expect.<sup>1</sup>

Looking at the results of varying dimension and reconstruction (changing one at a time or zeroing out) shows that each dimension can have a large effect. Most interestingly

is however that similar shapes are not near each other in the latent. We can see this because the reconstructions when varying a dimension jump from one shape, size, location to another, maybe even ghosting two shapes as it fades over 3.

### D. Autoencoder vs Variational autoencoder

Using the metrics discussed above we can measure the generative performance of both models quantitatively. The comparison is somewhat fair because both models train for about a minute and I spent similar amount of time hyperparameter tuning both models.

TABLE I  
COMPARISON OF GENERATIVE PERFORMANCE BETWEEN THE VAE AND AE MODELS. BOTH THE GENERATIVE PERFORMANCE AND RECONSTRUCTION PERFORMANCE ARE BETTER FOR THE VAE MODEL.

Metric	VAE	AE
Avg. Nearest Neighbor Distance	2.3478	13.5963
Coverage (%)	57.7%	0.0%
D_kl (bits)	8.1427	17.6513
Reconstruction Error (bits)	630.4751	74109.5781

TODO: Add in analysis of the table.

## IV. CONCLUSION

TODO: Add in conclusion.

<sup>1</sup>As a small experiment I tested out training the AE with only 4d. This resulted in some higher correlations between known features and latent dimensions, but still no clean separation of the 4 features.

## STATEMENT

TODO: add colab link The code and report of the Assignment was solely completed by myself (Thompson James). The complete source code can be found here [https://gitea.james-server.duckdns.org/james/AIML425\\_assignment\\_3](https://gitea.james-server.duckdns.org/james/AIML425_assignment_3), with a link to a colab notebook found here: TODO:addcolablink. A complete run through of the notebook takes about an hour to complete on a university lab machine.

I completed my work using the following tools:

- **Jupyter Notebook [1] and JupyterText [2]:** For interactive development and hosting.
- **LaTeX:** For writing the report.
- **VSCode [3]:** As IDE.
- **JAX [4] and Flax [5]:** For implementing the neural network and training logic.
- **Matplotlib [6] and Pandas [7]:** For data visualization and management.

## REFERENCES

- [1] T. Kluyver, B. Ragan-Kelley, F. Pérez, B. Granger, M. Bussonnier, J. Frederic, K. Kelley, J. Hamrick, J. Grout, S. Corlay, P. Ivanov, D. Avila, S. Abdalla, and C. Willing, “Jupyter Notebooks – a publishing format for reproducible computational workflows,” in *Positioning and Power in Academic Publishing: Players, Agents and Agendas*, F. Loizides and B. Schmidt, Eds. IOS Press, 2016, pp. 87–90.
- [2] M. Wouts, “Mwouts/jupyter,” Aug. 2025.
- [3] “Microsoft/vscode,” Microsoft, Jul. 2025.
- [4] J. Bradbury, R. Frostig, P. Hawkins, M. J. Johnson, C. Leary, D. Maclaurin, G. Necula, A. Paszke, J. VanderPlas, S. Wanderman-Milne, and Q. Zhang, “JAX: Composable transformations of Python+NumPy programs,” 2018.
- [5] J. Heek, A. Levskaya, A. Oliver, M. Ritter, B. Rondepierre, A. Steiner, and M. van Zee, “Flax: A neural network library and ecosystem for JAX,” 2024.
- [6] J. D. Hunter, “Matplotlib: A 2D graphics environment,” *Computing in Science & Engineering*, vol. 9, no. 3, pp. 90–95, 2007.
- [7] The pandas development team, “Pandas-dev/pandas: Pandas.”

## APPENDIX

TABLE II  
HYPERPARAMETERS USED FOR TRAINING THE VAE AND AE MODELS.

Hyperparameter	VAE	AE
Latent dimension	32	10
Encoder architecture	[1000, 1000, 1000, 500]	[2000, 2000, 2000]
Decoder architecture	[500, 1000, 1000, 1000]	[2000, 2000, 2000]
Optimizer	Adam(lr=0.001)	Adam(lr=0.001)
Minibatch size	512	512
Number of epochs	1000	200
Regularization weight	0.65	1
Dropout	Yes (p=0.1)	Yes (p=0.2)
Layer norm	No	No
Latent space noise	No	Yes ( $\sigma = 0.2$ )

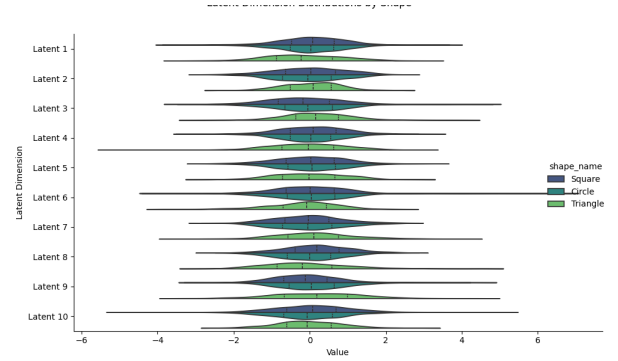


Fig. 1. Distribution of latent dimensions grouped by shape for the AE model. No strong correlation between shape and any latent dimension is visible.

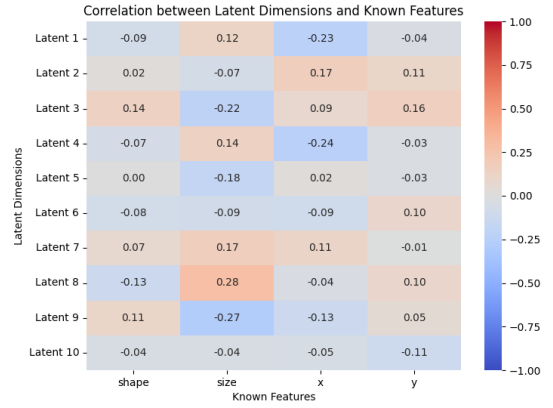


Fig. 2. Correlation between known features and latent dimensions for the AE model. Some strong correlation is visible, particularly for y position.

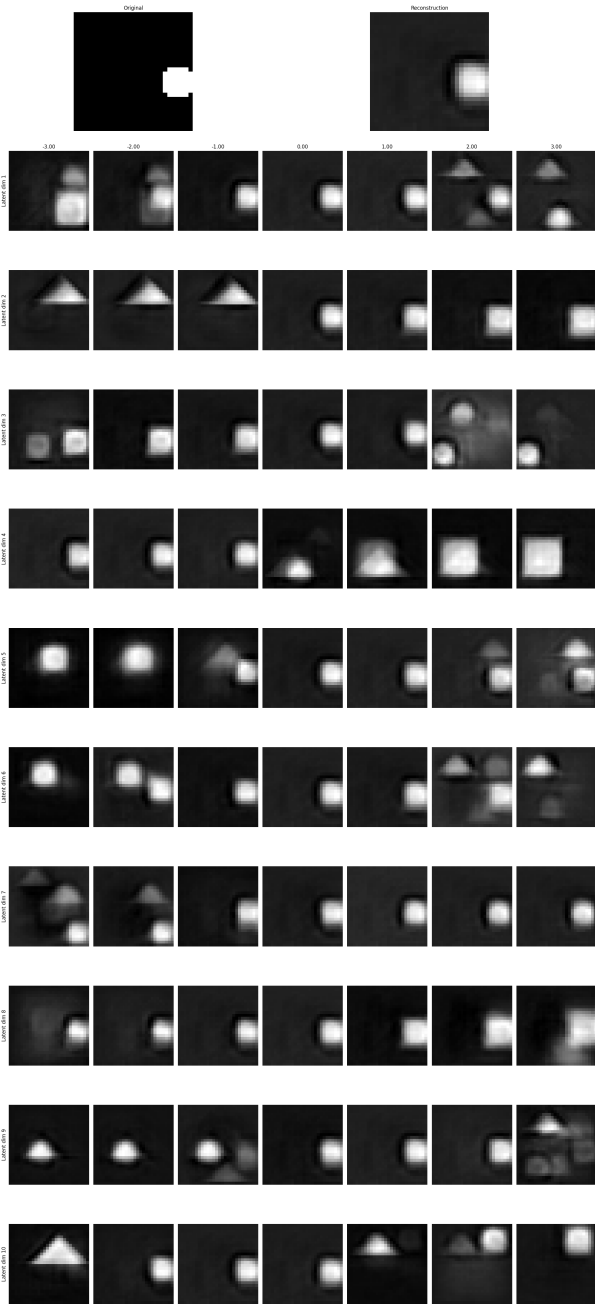


Fig. 3. Latent space traversals for the AE model. Each row shows the effect of varying one latent dimension from -3 to 3 while keeping the others constant.