

Studying data variability in variational autoencoders using a chain model

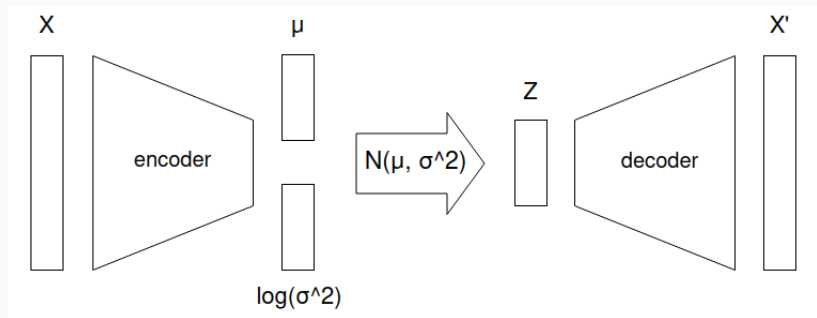
Tommaso Tarchi

January 16, 2024

University of Trieste

Introduction

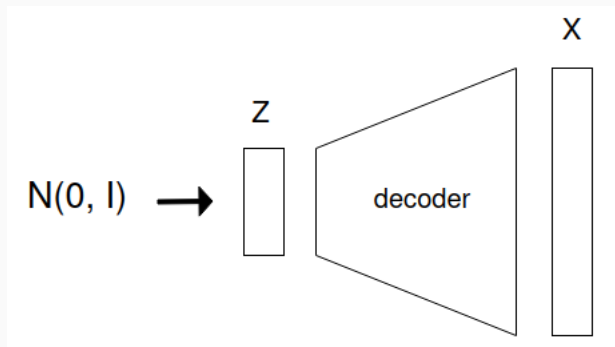
General architecture of a VAE



$$\text{LOSS}(\theta) = \text{MSE}_{\theta}(x, x_{\text{recon}}) + \text{KL}[q_{\theta}(z|x) || p(z)],$$

$$\text{with } p(z) = N(z|0, I)$$

Generative mode



What if...

But, what if we train a VAE on a dataset that was generated by another VAE?

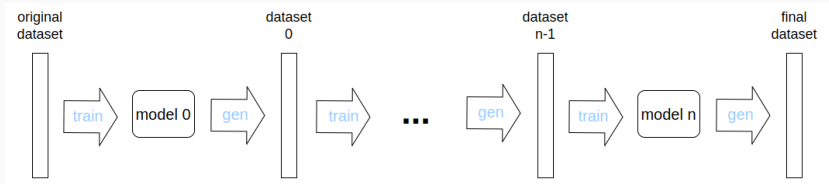
What if...

But, what if we train a VAE on a dataset that was generated by another VAE?

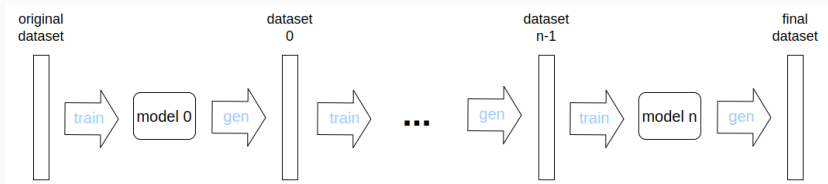
Possible questions:

- Is the new dataset still representative of the original dataset the first model was trained on?
- Are any bias introduced by this process?
- Will the new dataset have the same variability of the one it was trained on?

VAEs chain



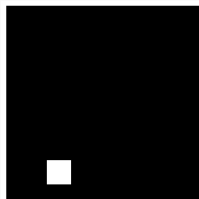
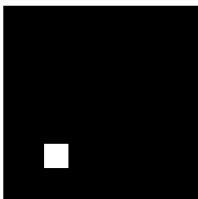
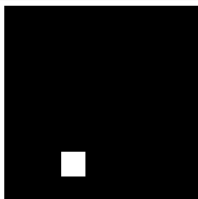
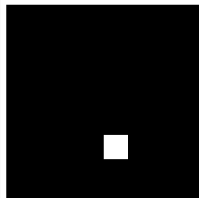
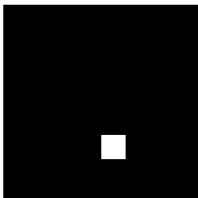
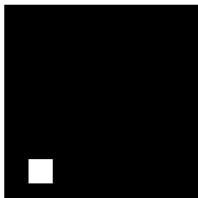
VAEs chain



Possible systematic effects introduced by VAEs will be amplified.

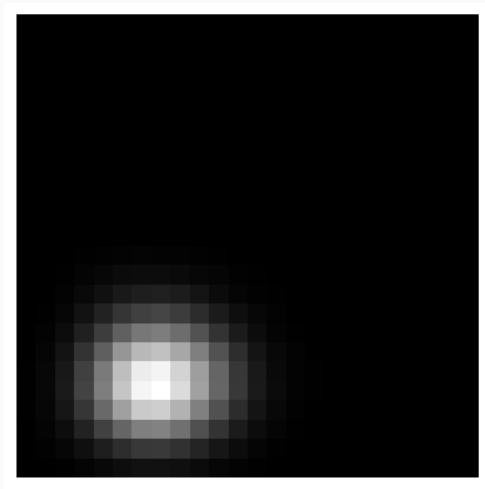
Our implementation

Initial dataset



24x24 grids, x-distribution: $\text{Binom}(24, 0.3)$, y-distribution: $\text{Binom}(24, 0.8)$

Dataset's distribution



24x24 grids, x-distribution: $\text{Binom}(24, 0.3)$, y-distribution: $\text{Binom}(24, 0.8)$

- One-layer fully connected ($FC1$)
- Two-layers fully connected ($FC2$)
- Convolutional with max-pooling ($Conv$)

How to evaluate results

- **Distribution difference:** difference between train dataset's and generated dataset's distributions (checks for biases)

How to evaluate results

- **Distribution difference:** difference between train dataset's and generated dataset's distributions (checks for biases)
- **Variability:** diversity of datasets (checks for variability reductions)

How to evaluate results

- **Distribution difference:** difference between train dataset's and generated dataset's distributions (checks for biases)
- **Variability:** diversity of datasets (checks for variability reductions)
- **Visual inspection:** checking random images visually (checks images to be of the right kind)

How to evaluate results

- **Distribution difference:** difference between train dataset's and generated dataset's distributions (checks for biases)
- **Variability:** diversity of datasets (checks for variability reductions)
- **Visual inspection:** checking random images visually (checks images to be of the right kind)

This method is NOT perfect!

Results

Single model

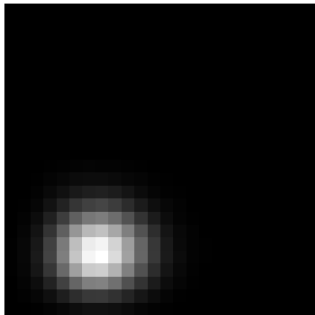


(a) Batch from original synthetic dataset

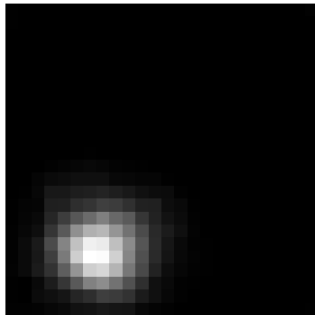


(b) Batch from dataset generated by VAE

Single model



(a) Distribution of original synthetic dataset



(b) Distribution of dataset generated by VAE

Whole chain



(a) Original dataset of the chain



(b) Dataset 1 of the chain



(c) Dataset 2 of the chain



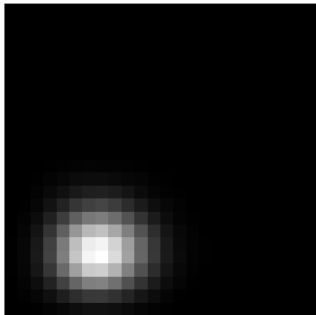
(d) Dataset 3 of the chain



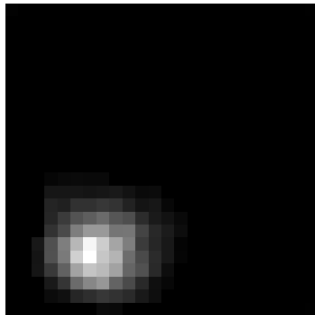
(e) Dataset 4 of the chain



(f) Final dataset of the chain

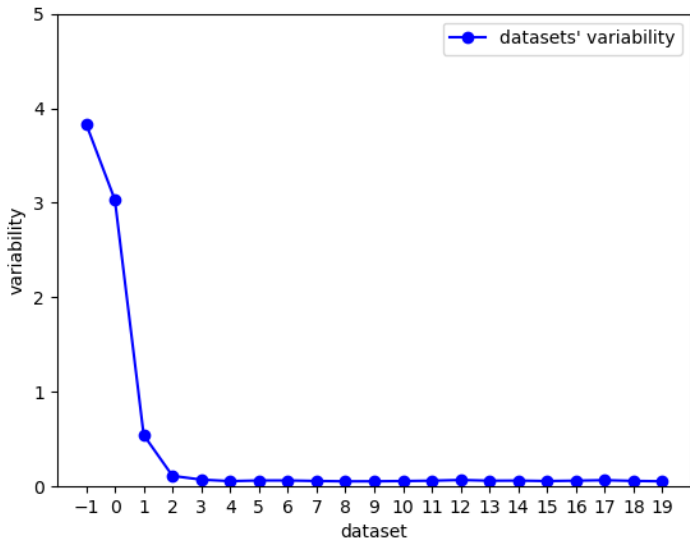


(a) Distribution of original synthetic dataset



(b) Distribution of dataset generated by last model in chain

Whole chain



- Images get lighter and more spread
- No evident bias is introduced in the data distribution (the data distribution is well "captured")
- **Variability is lost**

How can we fix this?

Possible causes

- **Regularization term:** the distribution in the latent space is forced to be close to a standard normal

Possible causes

- **Regularization term:** the distribution in the latent space is forced to be close to a standard normal
- **Dimensionality reduction:** less fundamental features are lost in latent space representation

Possible causes

- **Regularization term:** the distribution in the latent space is forced to be close to a standard normal
- **Dimensionality reduction:** less fundamental features are lost in latent space representation
- **Noisy reconstruction:** noise causes images to be more and more spread, may leading to more and more similar reconstructions

Possible causes

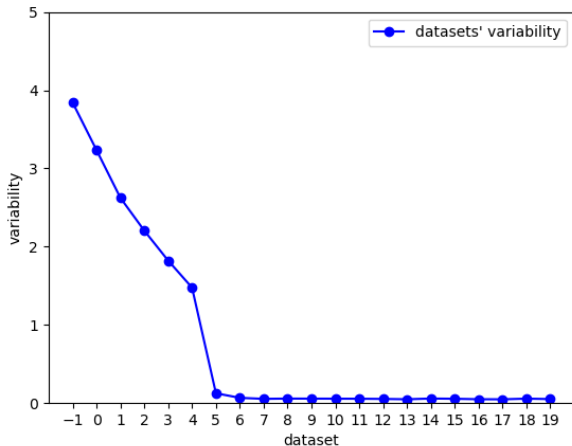
- **Regularization term:** the distribution in the latent space is forced to be close to a standard normal
- **Dimensionality reduction:** less fundamental features are lost in latent space representation
- **Noisy reconstruction:** noise causes images to be more and more spread, may leading to more and more similar reconstructions
- **Autoencoders' symmetry:** using mirror architectures for encoders and decoders may lead to unpredicted systematic patterns and behaviour

Possible causes

- **Regularization term:** the distribution in the latent space is forced to be close to a standard normal
- **Dimensionality reduction:** less fundamental features are lost in latent space representation
- **Noisy reconstruction:** noise causes images to be more and more spread, may leading to more and more similar reconstructions
- **Autoencoders' symmetry:** using mirror architectures for encoders and decoders may lead to unpredicted systematic patterns and behaviour
- **All (or some) of them combined**

$$\text{LOSS}(\theta) = \text{MSE}_{\theta}(x, x_{\text{recon}}) + \lambda \text{KL}[q_{\theta}(z|x) || p(z)]$$

Regularization constant tuning - results



Variability over models of the chain using $\lambda = 0.3$

Regularization constant tuning - results



(a) Batch from dataset generated by last model of the chain with $\lambda = 1$ (**untuned model**)



(b) Batch from dataset generated by last model of the chain with $\lambda = 0.3$ (**tuned model**)

Regularization constant tuning - conclusion

Pros:

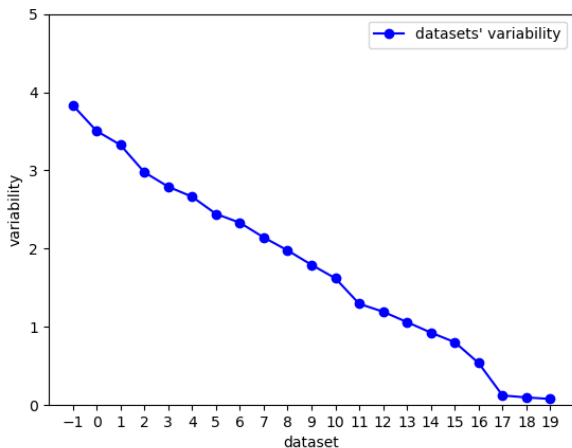
- helps maintaining variability longer
- easy to implement
- easy to tune ("linear" behaviour)

Cons:

- less precise outputs (latent space distribution less close to a standard normal)
- eventually leads to zero variability

$$\begin{aligned} \text{LOSS}_{batch}(\theta) = & \sum_i \text{MSE}_{\theta}(x_i, x_i^{recon}) + \sum_i \text{KL}[q_{\theta}(z_i|x_i) || p(z_i)] \\ & + K \left| \sum_{i,j} \text{MSE}(x_i, x_j) - \sum_{i,j} \text{MSE}_{\theta}(x_i^{recon}, x_j^{recon}) \right| \end{aligned}$$

Variability loss term - results



Dataset variability over models of the chain using $K = 10$

Variability loss term - results



(a) Batch from dataset generated by last model of the chain with $\lambda = 1$ and $K = 0$ (**no tuning**)



(b) Batch from dataset generated by last model of the chain with $\lambda = 0.3$ and $K = 0$ (**regularization tuning**)



(c) Batch from dataset generated by last model of the chain with $\lambda = 1$ and $K = 10$ (**variability loss tuning**)

Variability loss term - conclusion

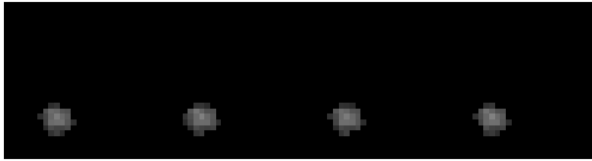
Pros:

- helps maintaining variability longer
- easy to tune
- **better variability-precision trade-off than regularization term tuning**
- **less light images**

Cons:

- still, eventually leads to zero variability

Denoising



(a) Non-denoised data



(b) Denoised data ($thres = 50\%$)

Improving non-optimal parameters with denoising

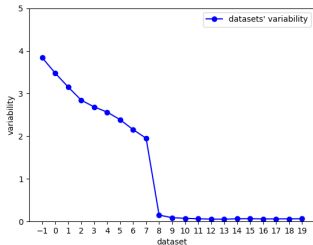


(a) Batch from fifth dataset without denoising ($\lambda = 0.7$ and $K = 10$)

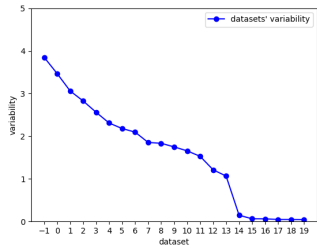


(b) Batch from fifth dataset with denoising ($\lambda = 0.7$ and $K = 10$, $\text{thres}=5\%$)

Improving non-optimal parameters with denoising

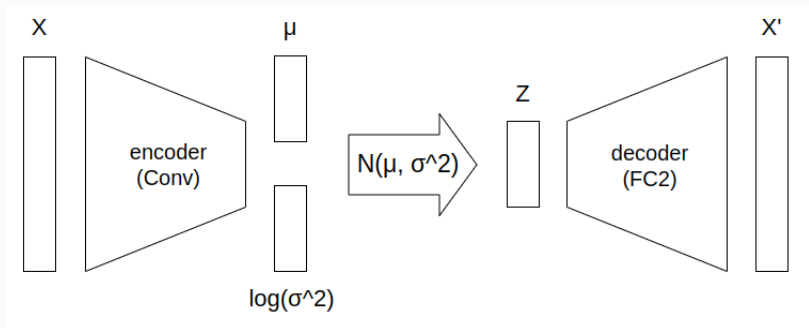


(a) Dataset variability without denoising ($\lambda = 0.7$, $K = 10$)

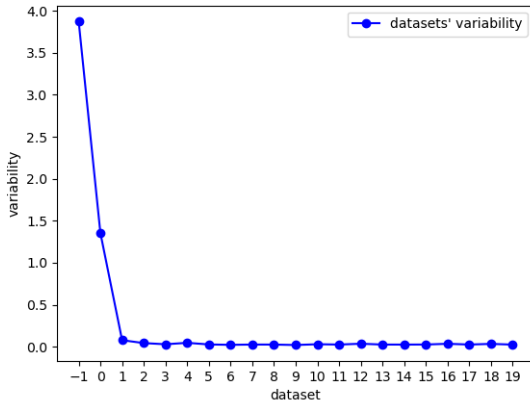


(b) Dataset variability with denoising ($\lambda = 0.7$, $K = 10$, $\text{thres}=5\%$)

Asymmetric architecture

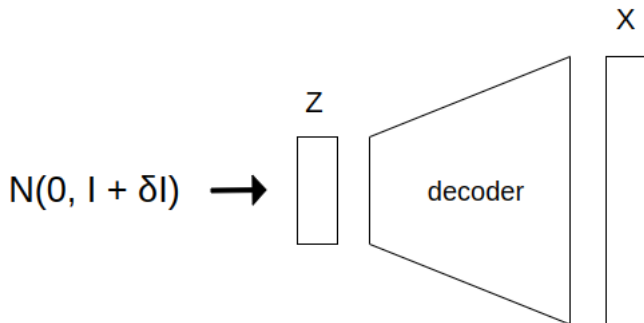


Asymmetric architecture - results

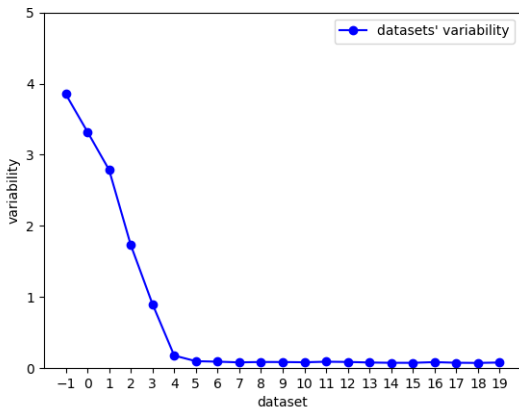


Dataset variability over models of the chain with asymmetric architecture

Wider distribution in generative mode

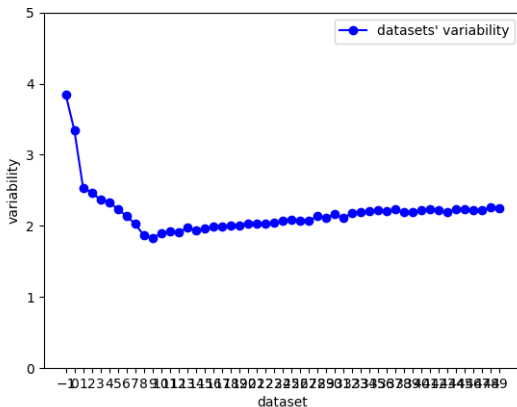


Wider distribution in generative mode - small increment



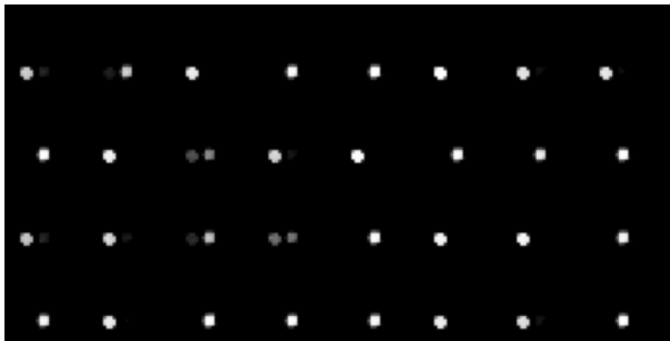
Dataset variability over models of the chain with $\delta = 0.3$, using denoising
(*thres* = 5%)

Wider distribution in generative mode - large increment



Dataset variability over models of the chain with $\delta = 0.7$, using denoising
(*thres* = 5%)

Wider distribution in generative mode - large increment



Batch from dataset generated by last model of the chain with $\delta = 0.7$,
using denoising (*thres* = 5%)

Wider distribution in generative mode - conclusion

Pros:

- helps maintaining variability longer
- **for some values of δ , keeps variability different from zero indefinitely**

Cons:

- **"classes formation"**
- may change distribution slightly
- very hard to tune (strongly non-linear behaviour)
- **unpredictable outcome on more complex datasets**
- may need denoising to work

Can we really fix this problem?

Can we really fix this problem?

Probably not! Due to information loss in latent space representation.

Can we really fix this problem?

Probably not! Due to information loss in latent space representation.

We can only either mitigate this effect or introduce "artificial" variability.

This is just a very simple example of a much larger and fundamental problem:

AI models trained on AI-generated data.

Thank you!