

## Future plans

My future plans is to run a hierarchical variational autoencoder on CelebA dataset as it is described below.

The great novelty of Variational Autoencoders was that they combined probabilistic models, i.e. variational Bayesian approach with deep learning techniques so that to perform efficient inference and continuous latent variables even in case of intractable posterior distributions, and large datasets.

With a reparameterization of the variational lower bound, a differentiable unbiased estimator of the lower bound can be achieved: the SGVB (Stochastic Gradient Variational Bayes) estimator can be used for efficient approximate posterior inference.

However, the interpretability of automatically discovered factorized representation of the independent data generative factors needed further modification on VAEs, i.e.: the introduction of an adjustable hyperparameter  $\beta$  that balances latent channel capacity and independence constraints with reconstruction accuracy.  $\beta$ -VAE with  $\beta > 1$  outperforms VAEs when the parameter was appropriately tuned [3.]. This modification limits the capacity of latent variable, which, combined with the pressure to maximize the log likelihood of the training data, should encourage the model to learn the most efficient representation of the data as the Kubleck-Leibler divergence term of the  $\beta$ -VAE objective function encourages conditional independence in the posterior: higher values of  $\beta$  should encourage learning a disentangled representation.

Quantitatively comparing the different unsupervised deep generative models a crucial point is the degree of disentanglement of the latent variables. In case of a disentangled representation one single latent units are sensitive to changes in single generative factors, while being relatively invariant to changes in other factors [4.] and knowledge about one factor can generalize to novel

configurations of other factors. [4.] also provides a protocol to quantitatively compare the degree of disentanglement learnt by different models.

Log likelihood of the data under the learnt model is a poor metric for evaluating disentangling in  $\beta$ -VAEs as latent channel capacity restriction ( $\beta > 1$ ) can lead to poorer reconstructions due to the loss of high frequency details when passing through a constrained latent bottleneck. [3.] propose a quantitative

metric that directly measures the degree of learnt disentanglement in the latent representation. It is important to note that a representation consisting of independent latents is not necessarily disentangled: i.e. PCA or ICA do not in general align with the data generative factors and hence may lack interpretability thus a simple cross-correlation calculation between the inferred latents would not suffice as a disentanglement metric.

Instead: inference is run on a number of images that are generated by fixing the value of one data generative factor while randomly sampling all others.

A low capacity linear classifier is used to identify the fixed factor and report the accuracy value as the final disentanglement metric score (as the independence and interpretability properties hold for the inferred representations, thus there will be less variance in the inferred latents that correspond to the fixed generative factor. Smaller variance in the latents corresponding to the target factor will make the job of this classifier easier, resulting in a higher score under the metric.

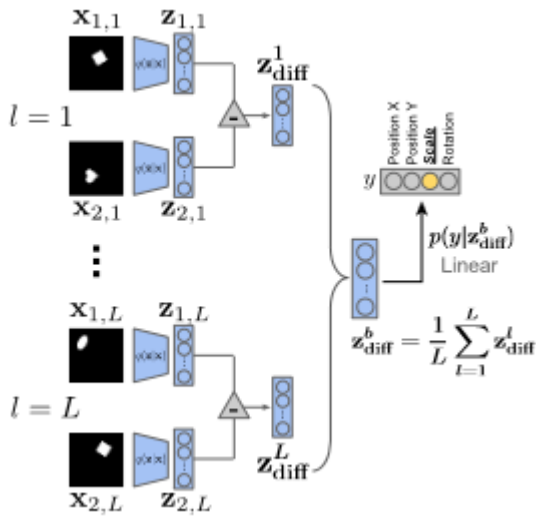


Figure 5: Schematic of the proposed disentanglement metric: over a batch of  $L$  samples, each pair of images has a fixed value for one target generative factor  $y$  (here  $y = scale$ ) and differs on all others. A linear classifier is then trained to identify the target factor using the average pairwise difference  $z_{diff}^b$  in the latent space over  $L$  samples.

### Tuning the $\beta$ coefficient

$\beta$  can be considered as a mixing coefficient for balancing the magnitudes of gradients from the reconstruction and the prior-matching components of the VAE lower bound formulation.  $\beta$  should be normalized by latent  $z$  size  $m$  and input  $x$  size  $n$  in order to compare its different values across different latent layer sizes and different datasets ( $\beta_{norm} = \beta * M/N$ ). We found that larger latent  $z$  layer sizes  $m$  requires higher constraint pressures (higher  $\beta$  values). Furthermore, the relationship of  $\beta$  for a given  $m$  is characterised by an inverted U curve. When  $\beta$

is too low or too high the model learns an entangled latent representation due to either too much or too little capacity in the latent  $z$  bottleneck:  $\beta > 1$  is necessary to achieve good disentanglement. However, if  $\beta$  is too high and the resulting capacity of the latent channel is lower than the number of data generative factors, then the learnt representation necessarily has to be entangled (as a low-rank projection of the true data generative factors will compress them in a non-factorial way to still capture the full data distribution well).

We proposed a principled way of choosing  $\beta$  for datasets with at least weak label information. If label information exists for at least a small subset of the independent data generative factors of variation, one can apply the disentanglement metric described in Sec. 3 to approximate the level of learnt disentanglement for various  $\beta$  choices during a hyperparameter sweep.

## Hierarchical VAEs

A more sophisticated approach for image generation is the nouveau VAE (NVAE) [5], a deep hierarchical VAE using depth-wise separable convolutions and batch normalization. NVAE is equipped with a residual parameterization of normal distributions and its training is stabilized by spectral regularization. NVAE uses depthwise convolutions in its generative model with (i) a new residual parameterization of the approximate posteriors. ii) stabilized training deep VAEs with spectral regularization, iii) practical solutions are used to reduce the memory burden of VAEs. iv) deep hierarchical VAEs can obtain state-of-the-art results on several image datasets and can produce high-quality samples even when trained with the original VAE objective.

The main building block of NVAE is depthwise convolutions that rapidly increase the receptive field of the network without dramatically increasing the number of parameters. (In depth-wise convolution, we use each filter channel only at one input channel. ) NVAE was the first successful application of VAEs to images as large as  $256 \times 256$  pixels.

[1.] Diederik P Kingma and Max Welling.  
Auto-encoding variational bayes.  
In The International Conference on Learning Representations (ICLR),  
2014.

<https://arxiv.org/pdf/1312.6114.pdf>

[2.] Diederik P Kingma and Max Welling.  
An introduction to variational autoencoders  
Foundations and Trends® in Machine Learning, 2019.

<https://arxiv.org/pdf/1906.02691.pdf>

[3.] Irina Higgins, Loic Matthey, Arka Pal, Christopher Burgess, Xavier Glorot,  
Matthew Botvinick, Shakir Mohamed, Alexander Lerchner  
ICLR 2017 conference submission, 2017.

Beta-vae: Learning basic visual concepts with a constrained variational  
framework

<https://openreview.net/forum?id=Sy2fzU9gl>

[4.] Y. Bengio, A. Courville, and P. Vincent.

Representation learning: A review and new perspectives. IEEE Transactions on  
Pattern Analysis & Machine Intelligence, 2013

[5.] Arash Vahdat, Jan Kautz

NVAE: A Deep Hierarchical Variational Autoencoder

<https://arxiv.org/pdf/2007.03898.pdf>

[6.] Tero Karras, Timo Aila, Samuli Laine, Jaakko Lehtinen Progressive Growing  
of GANs for Improved Quality, Stability, and Variation

<https://arxiv.org/pdf/1710.10196v3.pdf>

[7.] SZEGEDY, Christian, et al. Going deeper with convolutions.

In: Proceedings of the IEEE conference on computer vision and pattern  
recognition. 2015. p. 1-9.

[Going Deeper With Convolutions \(googleusercontent.com\)](https://arxiv.org/pdf/1502.00889v1.pdf)