

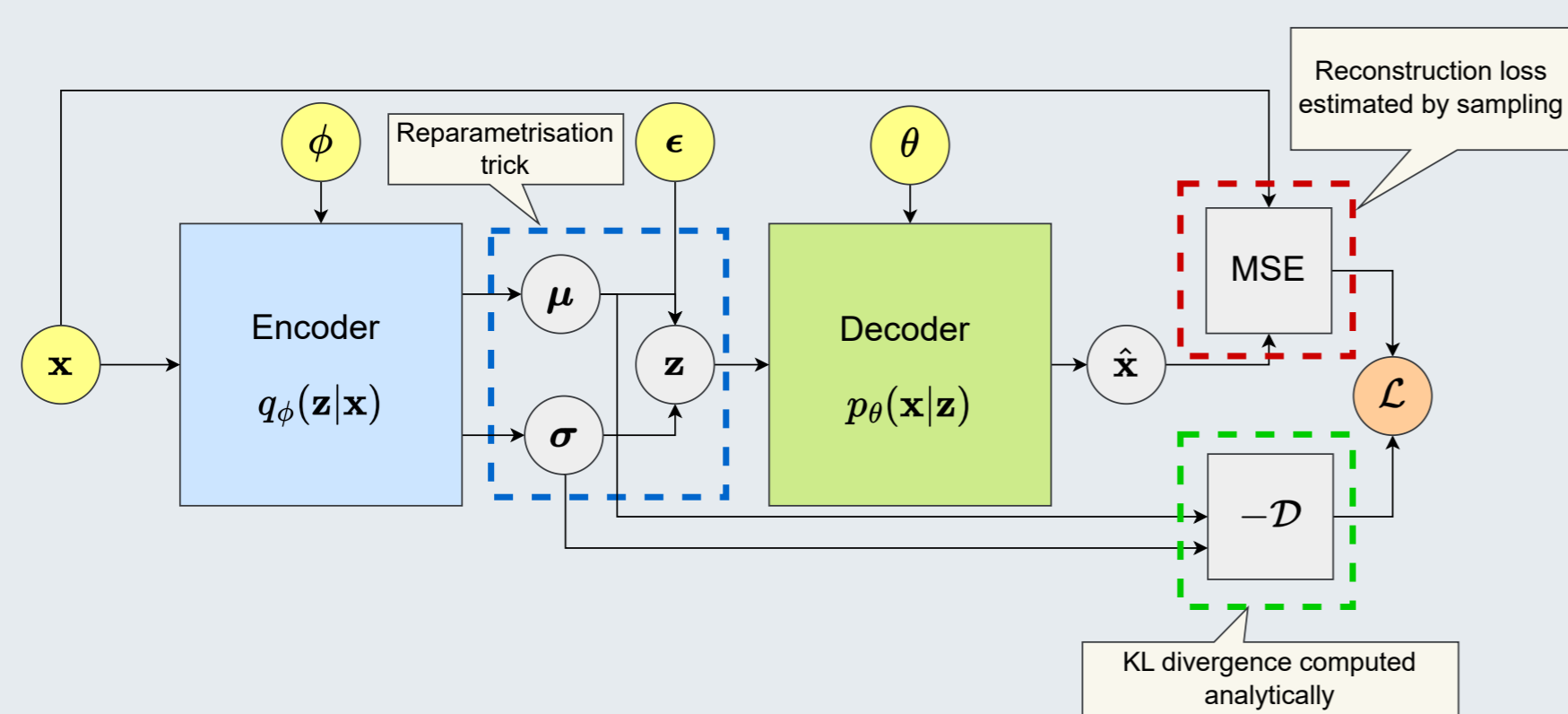
Lisa Bonheme
 Marek Grzes

Summary

Locatello et al. [3] observed a lower disentanglement in mean than sampled representations of Variational Autoencoders (VAEs). In this paper we:

- Analyse the problem through the lens of the polarised regime
- Show that the lower disentanglement of mean representations is due to (uninformative) passive variables
- Provide some recommendation for using mean representations on downstream tasks

What are Variational Autoencoders?

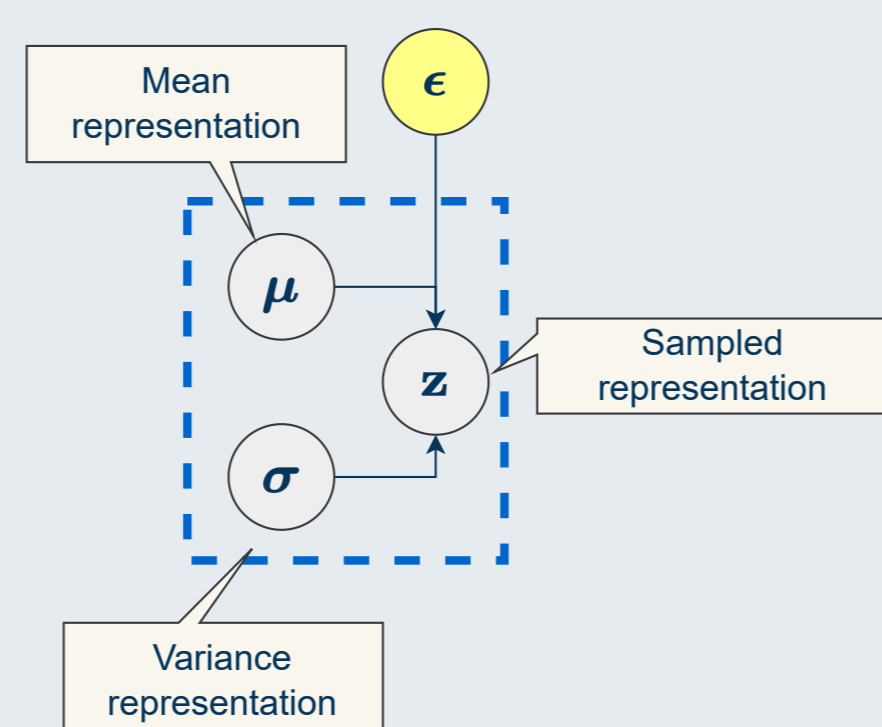


$$\mathcal{L}(\theta, \phi; \mathbf{x}) = \underbrace{\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})]}_{\text{reconstruction term}} - \underbrace{D_{\text{KL}}(q_{\phi}(\mathbf{z}|\mathbf{x}) \parallel p_{\theta}(\mathbf{z}))}_{\text{regularisation term}}.$$

The samples from the learned latent representation are obtained using the reparameterisation trick [2] such that $\mathbf{z} = \boldsymbol{\mu} + \boldsymbol{\sigma}^{1/2}\boldsymbol{\epsilon}$ where $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$.

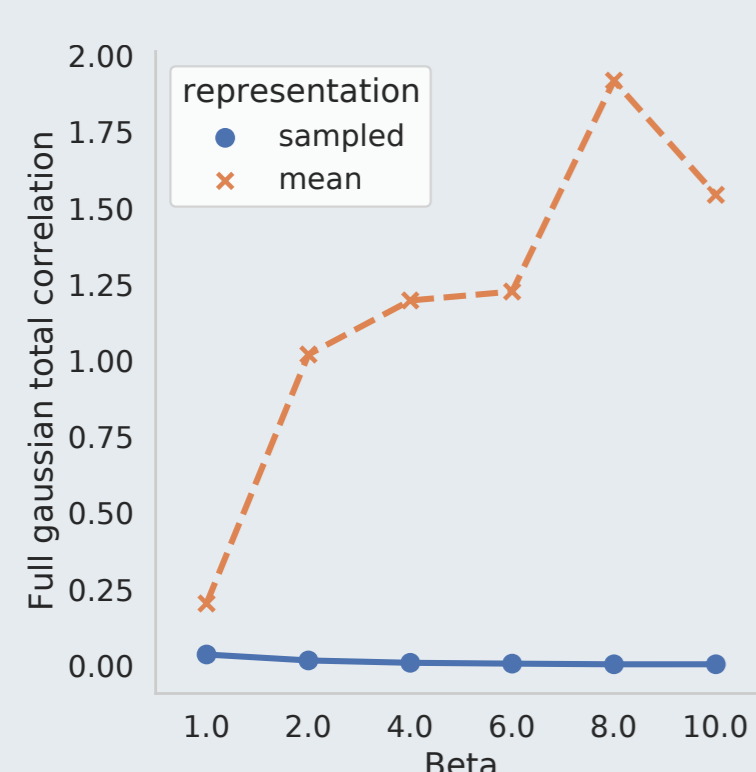
The polarised regime [1, 4]

- **Passive variables**
 $\mu_i \approx 0$, $\sigma_i \approx 1$, and $\mathbf{z}_i \sim \mathcal{N}(0, 1)$.
- **Active variables**
 $\sigma_i \approx 0$ and $\mathbf{z}_i \approx \mu_i$.

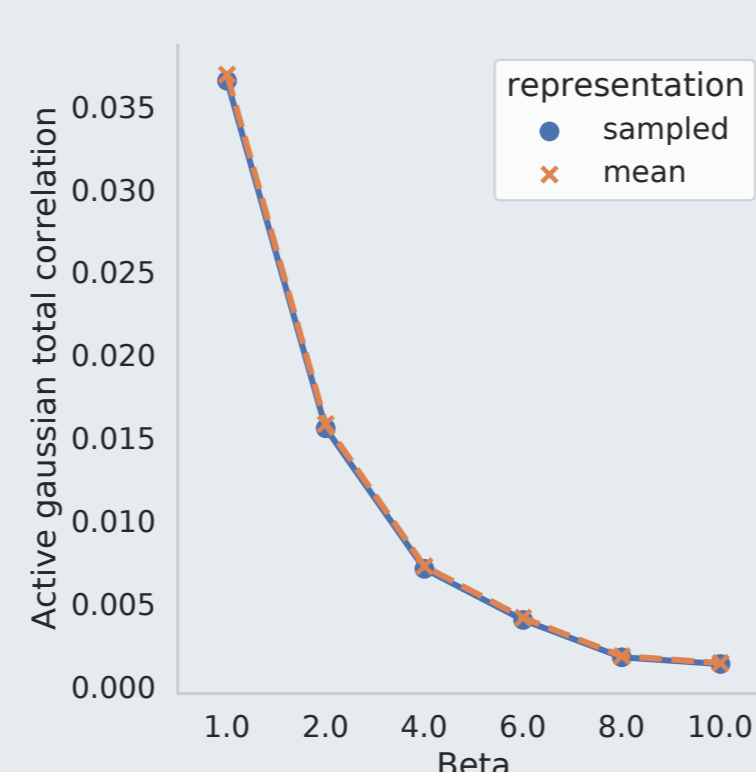


Truncation experiment

- Total correlation (TC) and averaged mutual information (MI) are higher in $\boldsymbol{\mu}$ than \mathbf{z} [3].
- For active variables $\mathbf{z}_i \approx \mu_i$.
- Discrepancies between $\boldsymbol{\mu}$ and \mathbf{z} must come from passive variables.



(a) Full representation

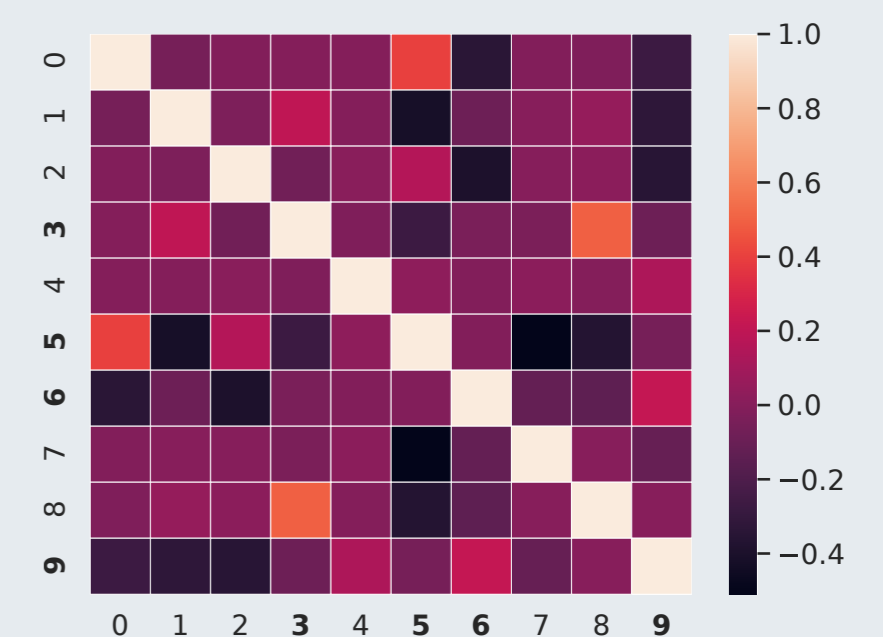


(b) Only active variables

Figure: TC of β -TC VAE on noisy dSprites

Passive variables correlation

- Passive variables are more correlated in $\boldsymbol{\mu}$ than \mathbf{z} .
- They should be removed before using $\boldsymbol{\mu}$ on downstream tasks.



Where does this correlation come from?

- Are variables passive because they are correlated?
- Does the correlation occur because the variables are passive?

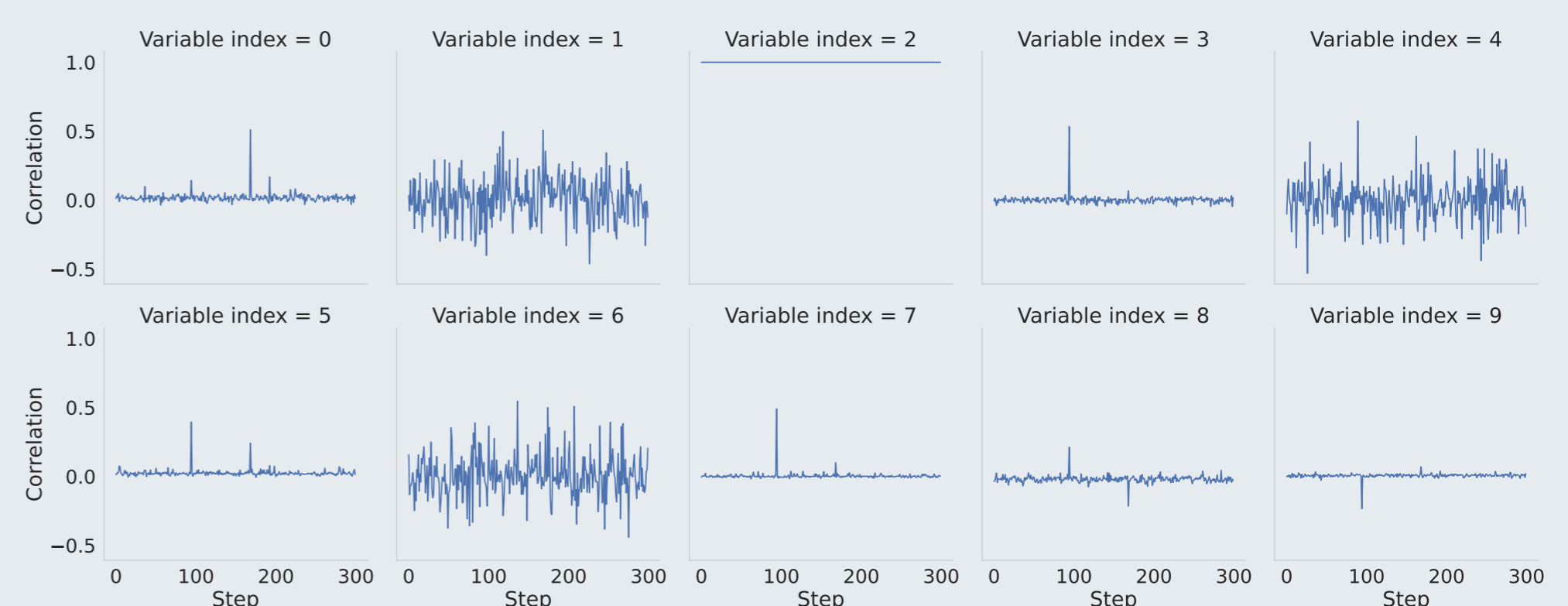


Figure: The correlation scores of the active variable at index 2 of the mean representation with all the other variables during the 300K training steps of a β -VAE with $\beta = 8$ trained on dSprites. We can see an increased correlation with all the passive variables (indexes 1, 4, and 6).

Conclusion

- Active variables are as disentangled in mean as in sampled representations
- Passive variables are highly correlated with various active variables
- An in-depth study of the learning dynamics of VAEs would be needed to explain this phenomenon
- Passive variables should be removed from mean representations before downstream tasks

More about the paper



References

- [1] B. Dai, Y. Wang, J. Aston, G. Hua, and D. Wipf. Connections with Robust PCA and the Role of Emergent Sparsity in Variational Autoencoder Models. *Journal of Machine Learning Research*, 19(41):1–42, 2018.
- [2] D. P. Kingma and M. Welling. Auto-Encoding Variational Bayes. In *International Conference on Learning Representations*, volume 2, 2014.
- [3] F. Locatello, S. Bauer, M. Lucic, G. Raetsch, S. Gelly, B. Schölkopf, and O. Bachem. Challenging Common Assumptions in the Unsupervised Learning of Disentangled Representations. In *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, 2019.
- [4] M. Rolínek, D. Zietlow, and G. Martius. Variational Autoencoders Pursue PCA Directions (by Accident). In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.