

Task 4: Variational Autoencoder (VAE)

Variational Auto Encoders learn probabilistic gaussian latent distribution of data. The encoder maps images to the latent space. The decoder reconstructs images from the latent distribution. Learning is made possible through the use of the reparameterization trick. I have explored different techniques to reconstruct the best representation of images. This includes exploring active z-dimensions, role of KL and reconstruction loss on regenerations and sampling from different distributions i.e. normal or laplacian.

1. Train the VAE

I have trained a VAE on Fashion-MNIST with a Gaussian decoder (MSE) and a standard Normal prior. Loss curves show stable optimization: reconstruction error fell from 28.63 to 14.88 MSE per sample while the KL term rose from 6.13 and plateaued at approx. 7.98 nats, yielding a total ELBO around 22.86 as in Fig. (4.1)

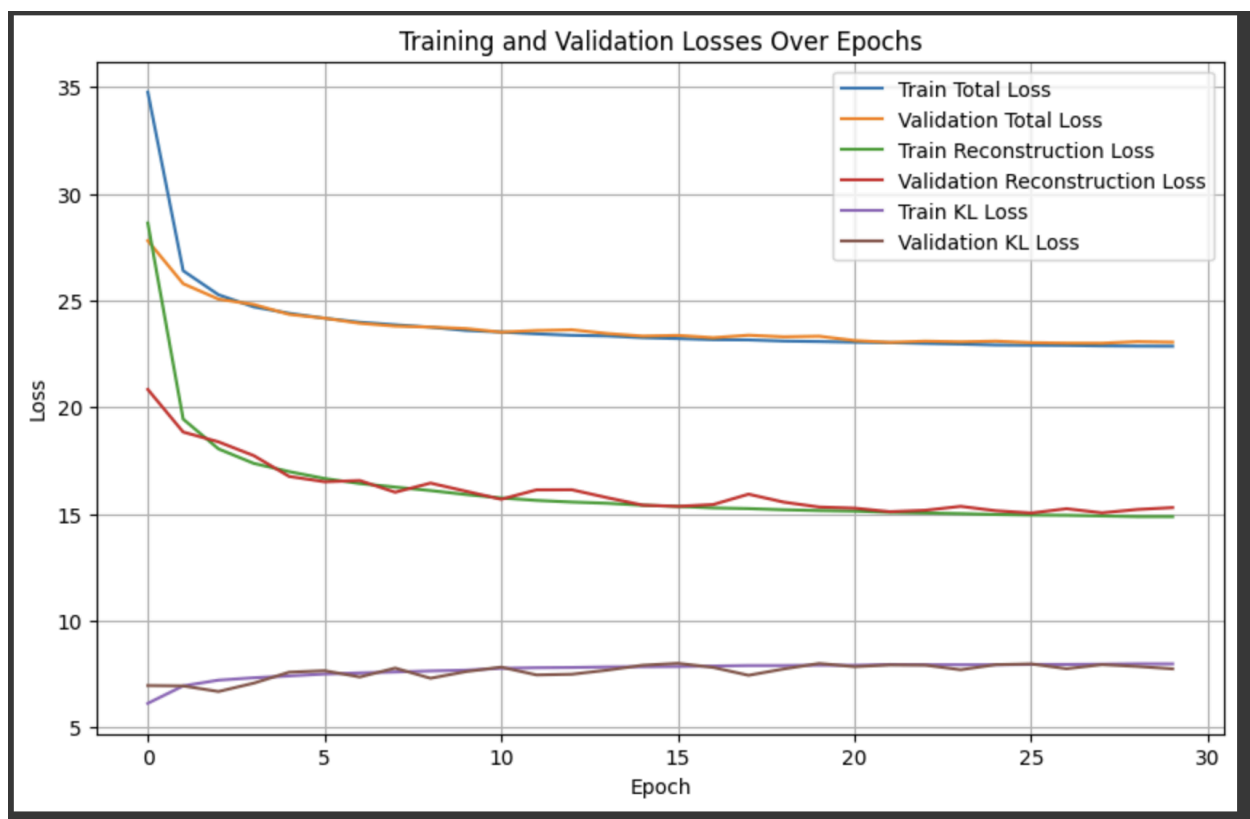


Fig (4.1)

2. Visualize Reconstructions and Generations

Reconstructions are compared with originals and it was observed that reconstructions are blurry and only preserves edges and silhouettes. We can see that textures and texts are all gone as Gaussian approximates pixels.

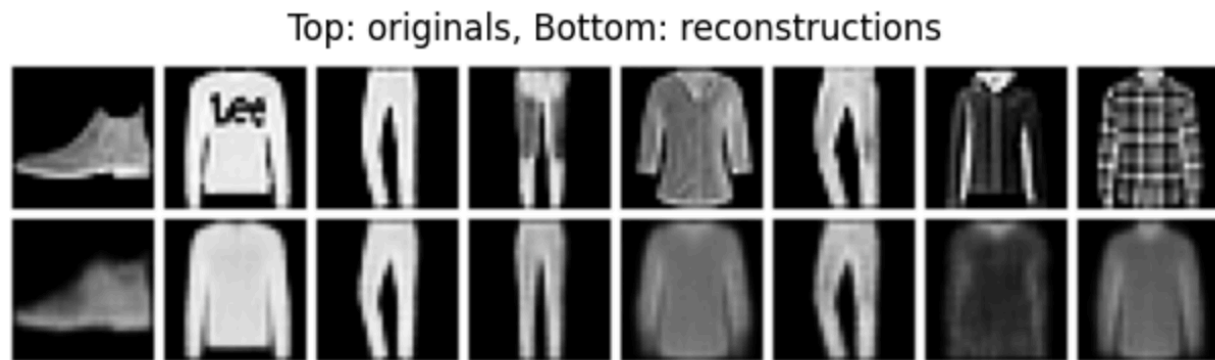


Fig (4.2)

3. Gaussian vs laplacian generations

Sampling from the Normal prior has produced low quality images, the reason could be a small number of epochs and because most dimensions were unused and only 6 out of 20 dimensions have contributed to the results. The encoder did not have enough to learn, sampling from a Laplace prior also degrades quality due to prior-decoder mismatch. Since images are reconstructed from Gaussian latent space of 20dim, the images are blurry and edges aren't refined.

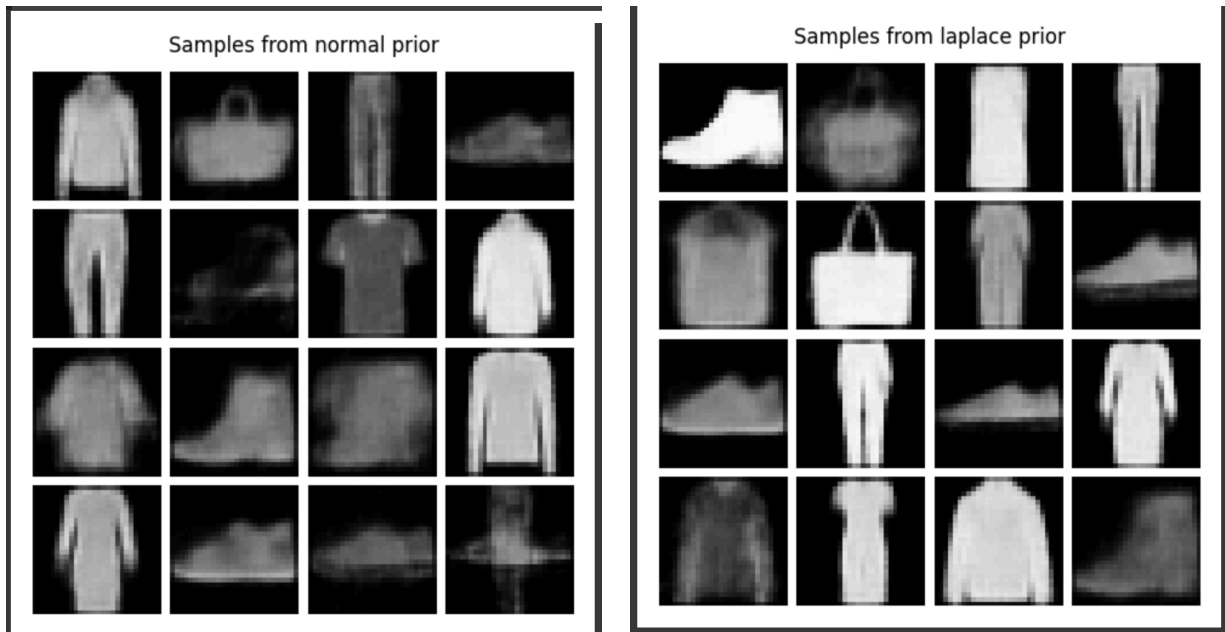
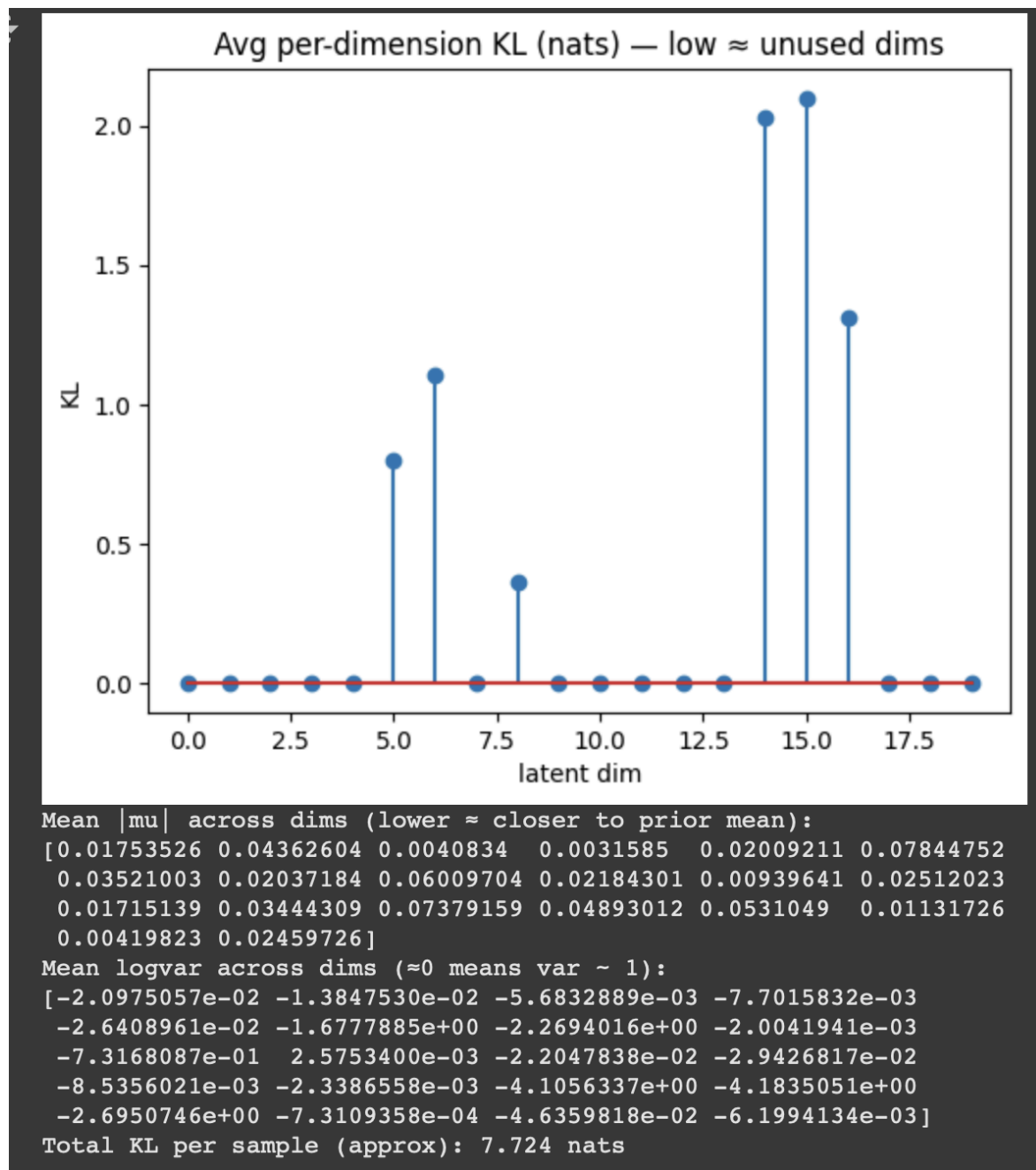


Fig (4.3)

4. Posterior Collapse

Plot of Fig (4.2) clearly shows which latent coordinates are used by model. For most dimensions KL is zero and only 6 dimensions show a spike with considerable KL values. The stats show that mean per dimension is small and log per variance is highly negative for the active dimensions. From the stats we can see that KL is largely variance-driven.



Fig(4.4)

4. Mitigating Posterior Collapse

To mitigate the effect of posterior collapse techniques of warm-up = 10 and free bits = 0.05 are introduced. This was to give more room for the encoder to learn and map images in latent space. Resultantly images reconstructed have better quality and instead of 4 dimensions now all 20 dimensions are used.

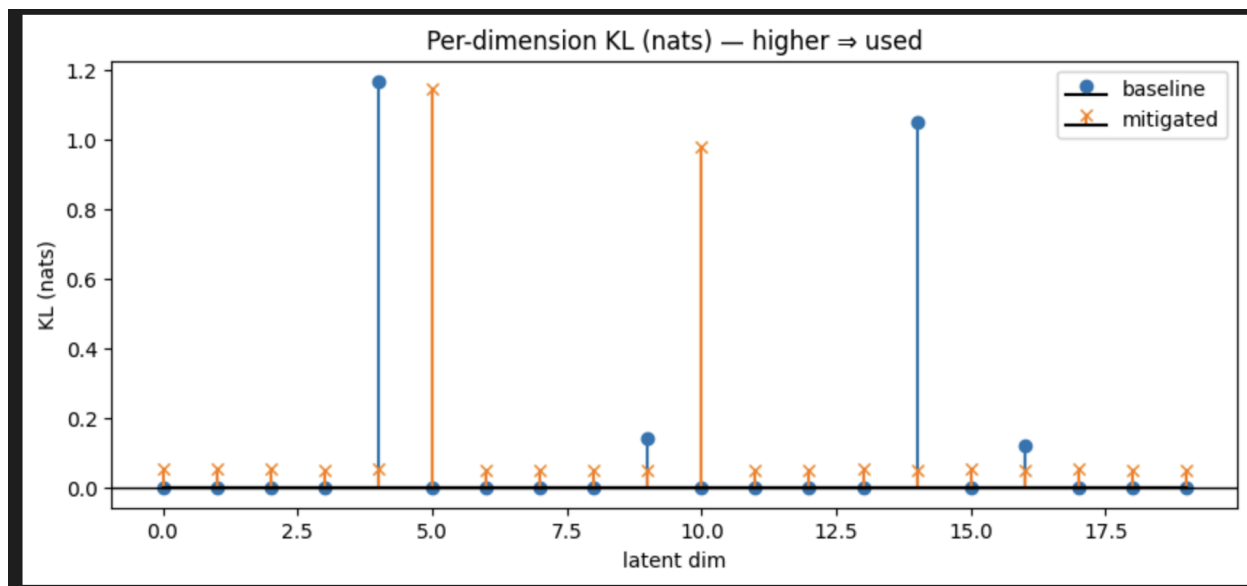
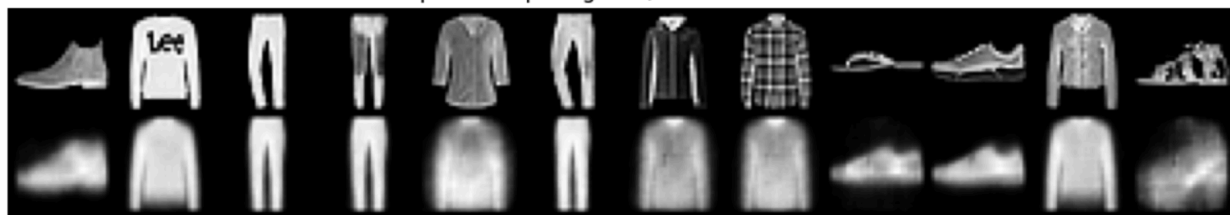
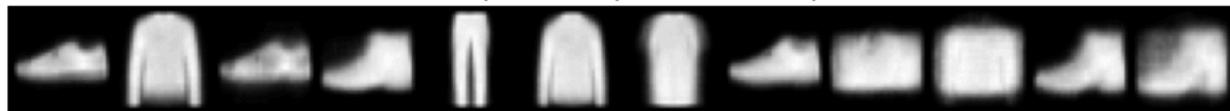


Fig (4.5)

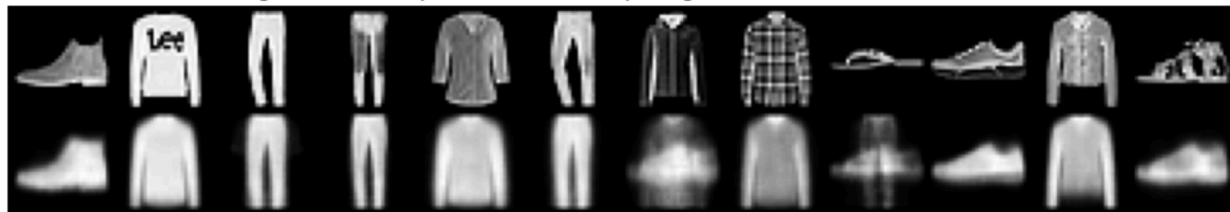
Baseline $\beta=1$ — Top: originals, Bottom: reconstructions



Baseline $\beta=1$ — Samples from Normal prior



Mitigated (warm-up + free-bits) — Top: originals, Bottom: reconstructions



Mitigated (warm-up + free-bits) — Samples from Normal prior

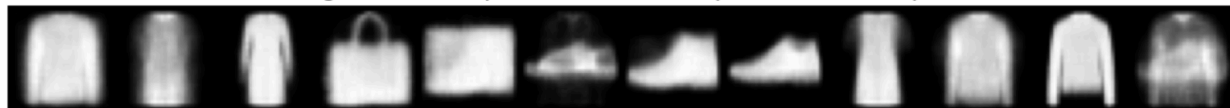


Fig (4.6)

Conclusion:

Per-dimension analysis reveals partial posterior collapse, despite a 20-D latent, only 5–6 dimensions carry substantial KL, largely via variance shrinkage (means near zero, strongly negative log-variances on active dims), therefore reconstructions preserve category and structure but lose fine details. WarmUp and free bits were introduced to enhance the quality of reconstructed images. Reconstruction is consistent with the information budget and the Gaussian objective. Overall, the model learns a compact, smooth latent representation that supports plausible generations and reconstructions.

Some other possible ways to improve the model can be:

- 1) Change of objective function, we can try using BCE with logits on FashionMNIST dataset as Bernoulli can preserve edges and details better than MSE
- 2) Increasing dimensions of latent space. It may increase KL but can reduce reconstruction error